

CS540 Introduction to Artificial Intelligence

Lecture 2

Young Wu

Based on lecture slides by Jerry Zhu, Yingyu Liang, and Charles Dyer

May 10, 2020

Feedback

Admin

- Please give me feedback on lectures, homework, exams on Socrative, room CS540.
- Please report bugs in homework, lecture examples and quizzes on Piazza.
- Please do NOT leave comments on YouTube.
- Email me (Young Wu) for personal issues.
- Email department chair (Dieter van Melkebeek) for issues with me.

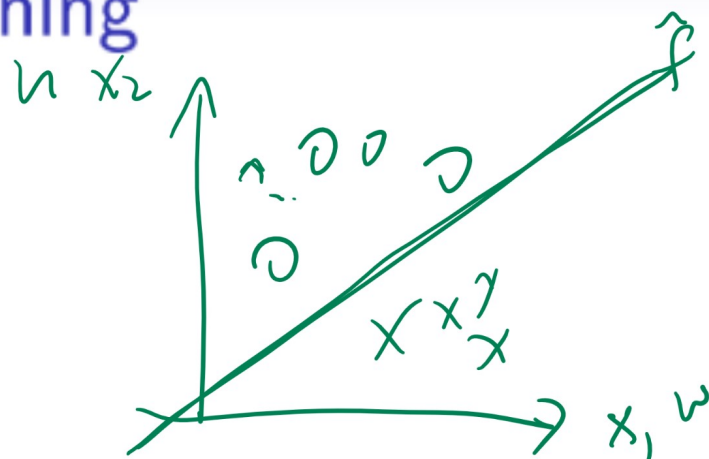
Feedback from Last Year

Admin

- Too much math.
- Time spent on math.
- Cannot understand my handwriting.
- Mistake on slides.
- More examples.

Supervised Learning

Motivation

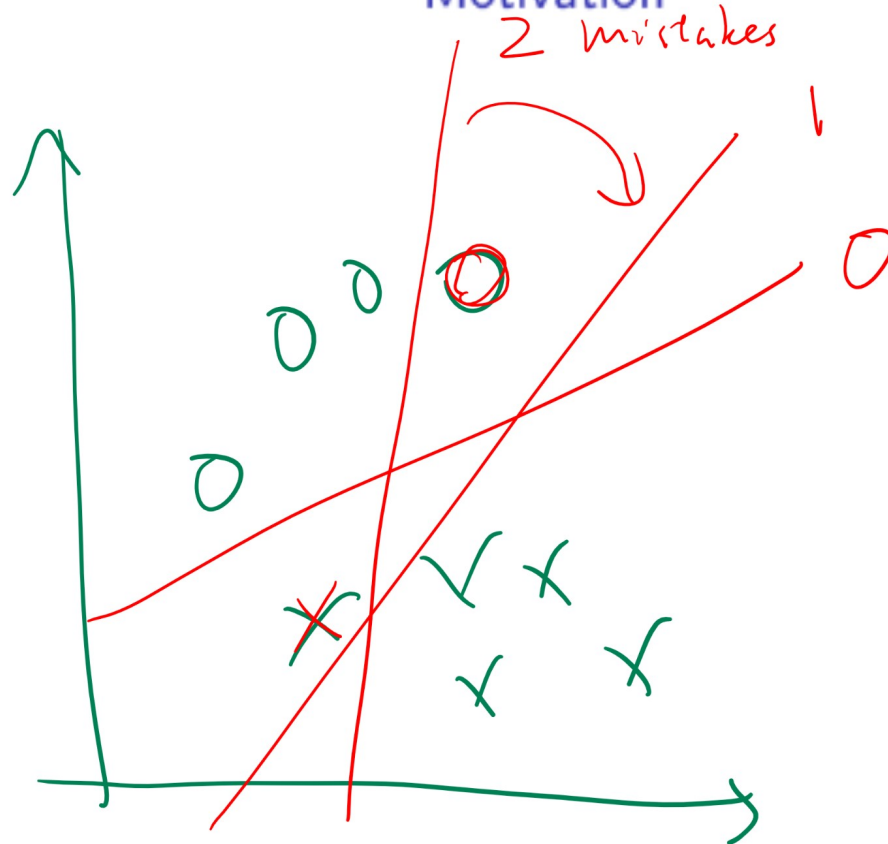


Data	Features (Input)	Output	-
Training	$\{(x_{i1}, \dots, x_{im})\}_{i=1}^{n'}$	$\{y_i\}_{i=1}^{n'}$	find "best" \hat{f}
	observable	known	-
Test	(x'_1, \dots, x'_m)	y'	guess $\hat{y} = \hat{f}(x')$
-	observable	unknown	-



Loss Function Diagram

Motivation



min #mistakes
f

Zero-One Loss Function

Motivation

minimize \rightarrow # mistakes
 argmin \Rightarrow f that makes
 min # mistakes

- An objective function is needed to select the "best" \hat{f} . An example is the zero-one loss.

$$\hat{f} = \arg \min_f \sum_{i=1}^n \mathbb{1}_{\{f(x_i) \neq y_i\}}$$

training samples (circled around n)
 # mistakes (circled around the indicator function)
 1 mistake \downarrow
 1 $\hat{y} \neq y$
 0 $\hat{y} = y$

- $\arg \min_f$ objective (f) outputs the function that minimizes the objective.
- The objective function is called the cost function (or the loss function), and the objective is to minimize the cost.

Squared Loss Function

Motivation

- Zero-one loss counts the number of mistakes made by the classifier. The best classifier is the one that makes the fewest mistakes.
- Another example is the squared distance between the predicted and the actual y value:

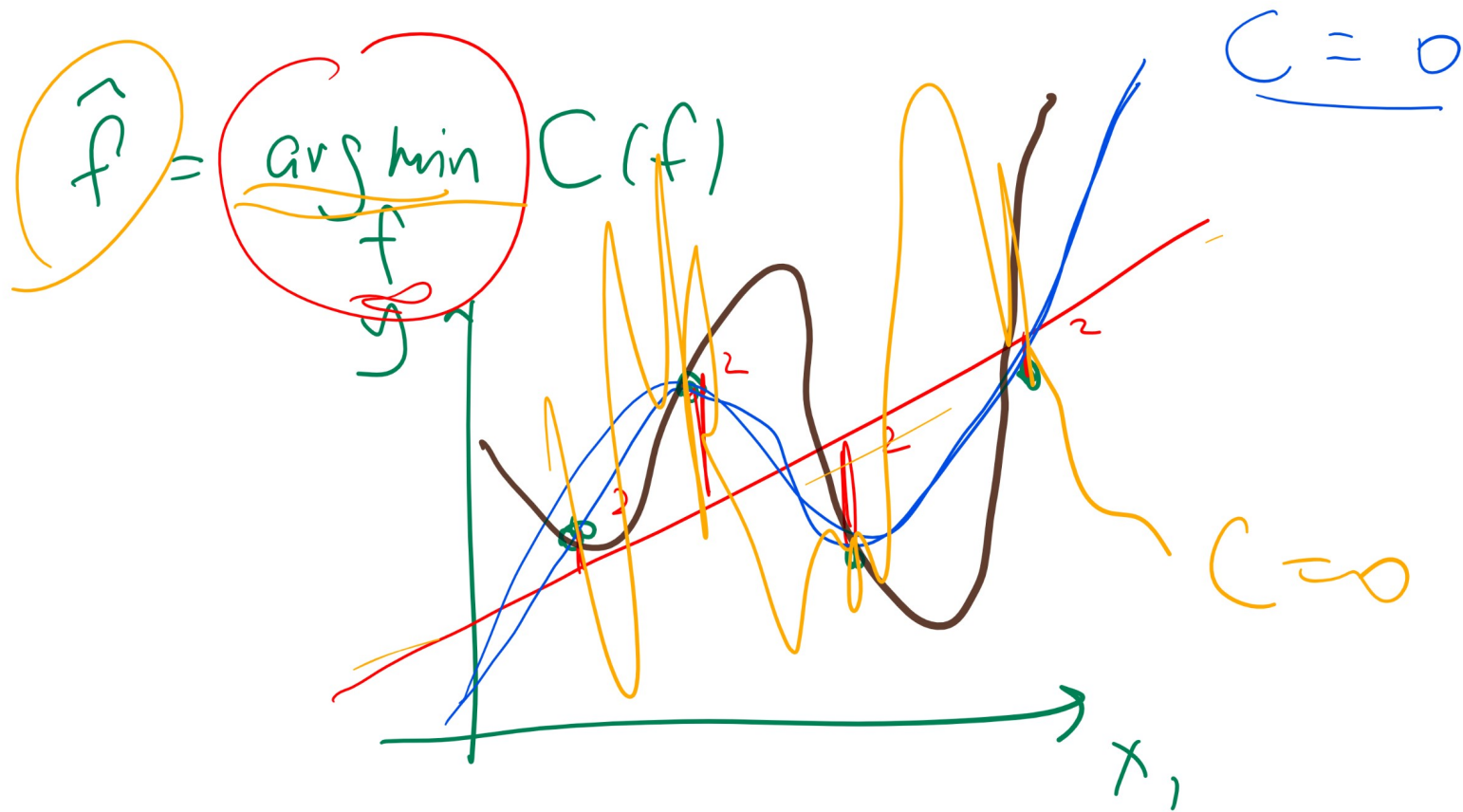
$$\hat{f} = \arg \min_f \frac{1}{2} \sum_{i=1}^n (f(x_i) - y_i)^2$$

Handwritten annotations illustrating the squared loss function:

- Green arrows point from the predicted values 0.6 and 0.5 to the squared differences 0.4 and 0.25 .
- A bracket under the sum indicates the total squared error.
- A note $\{ \hat{y}_i \neq y_i \}$ with arrows pointing to the predicted values 0.6 and 0.5 indicates that these predictions are incorrect.
- Vertical arrows on the right side of the equation point to the predicted values 0.6 and 0.5 .

Loss Functions Equivalence

Quiz



Hypothesis Space

Motivation

- There are too many functions to choose from.
- There should be a smaller set of functions to choose \hat{f} from.

$$\hat{f} = \arg \min_{f \in \mathcal{H}} \frac{1}{2} \sum_{i=1}^n (f(x_i) - y_i)^2$$

- The set \mathcal{H} is called the hypothesis space.

Linear Regression

Motivation

- For example, \mathcal{H} can be the set of linear functions. Then the problem can be rewritten in terms of the weights.

$$\left(\hat{w}_1, \dots, \hat{w}_m, \hat{b} \right) = \arg \min_{w_1, \dots, w_m, b} \frac{1}{2} \sum_{i=1}^n (a_i - y_i)^2$$

where $a_i = w_1 x_{i1} + w_2 x_{i2} + \dots + w_m x_{im} + b$

Handwritten notes: β above $\hat{w}_1, \dots, \hat{w}_m, \hat{b}$; \hat{f} below $\hat{w}_1, \dots, \hat{w}_m, \hat{b}$; w_1, \dots, w_m, b circled in the denominator of the sum.

- The problem is called (least squares) linear regression.

$$y = 0$$

Handwritten notes: $\frac{0}{1} \rightarrow \text{cat}$, $\frac{1}{2} \rightarrow \text{dog}$, $\frac{2}{3} \rightarrow \text{fox}$

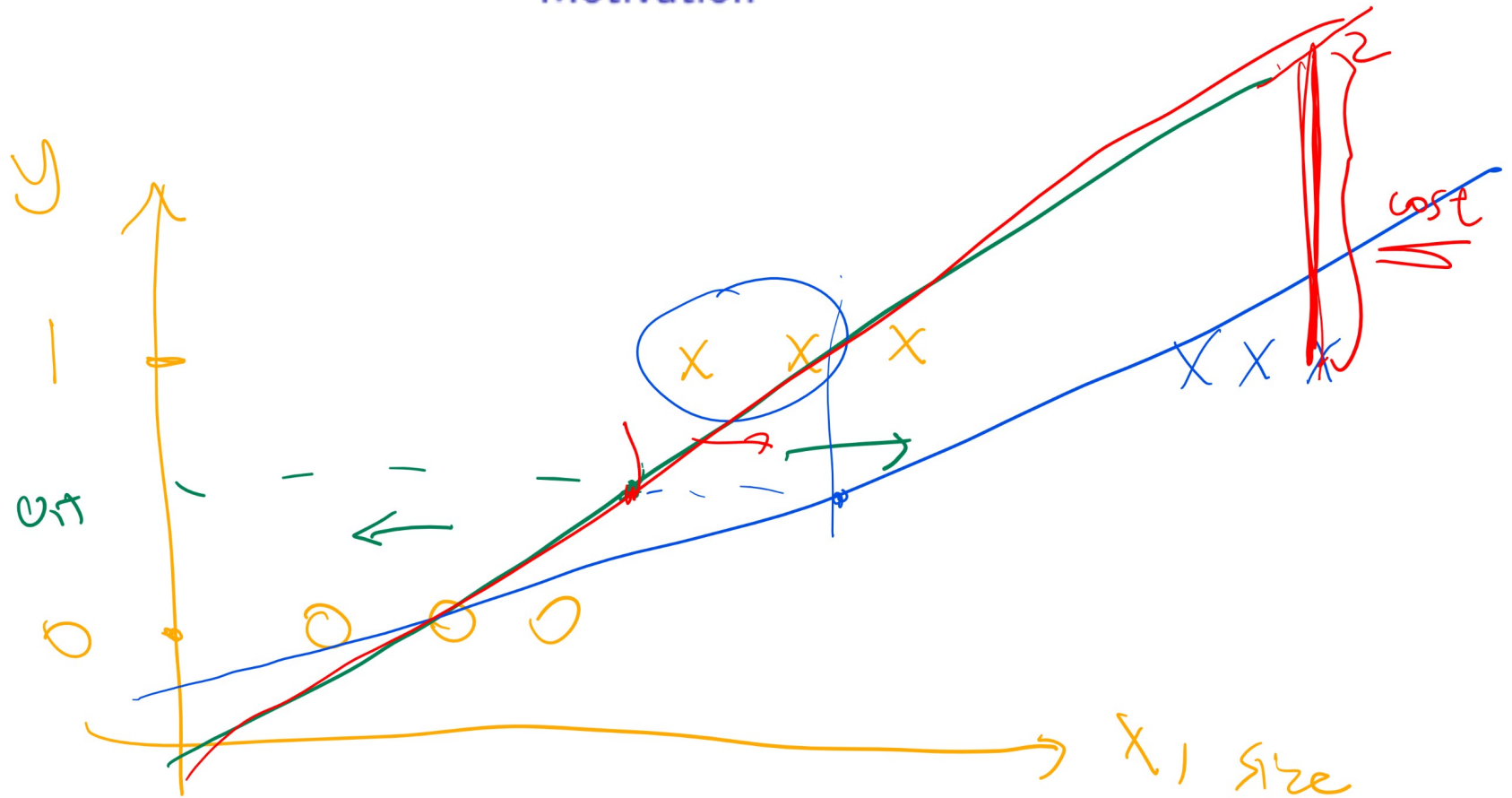
Binary Classification

Motivation

- If the problem is binary classification, y is either 0 or 1, and linear regression is not a great choice.
- This is because if the prediction is either too large or too small, the prediction is correct, but the cost is large.

Binary Classification Linear Regression Diagram

Motivation



Linear Threshold Unit

Motivation

- One simple choice is to use the step function as the activation function:

$$g(\boxed{\cdot}) = \mathbb{1}_{\{\boxed{\cdot} \geq 0\}} = \begin{cases} 1 & \text{if } \boxed{\cdot} \geq 0 \\ 0 & \text{if } \boxed{\cdot} < 0 \end{cases}$$

- This activation function is called linear threshold unit (LTU).

Sigmoid Activation Function

Motivation

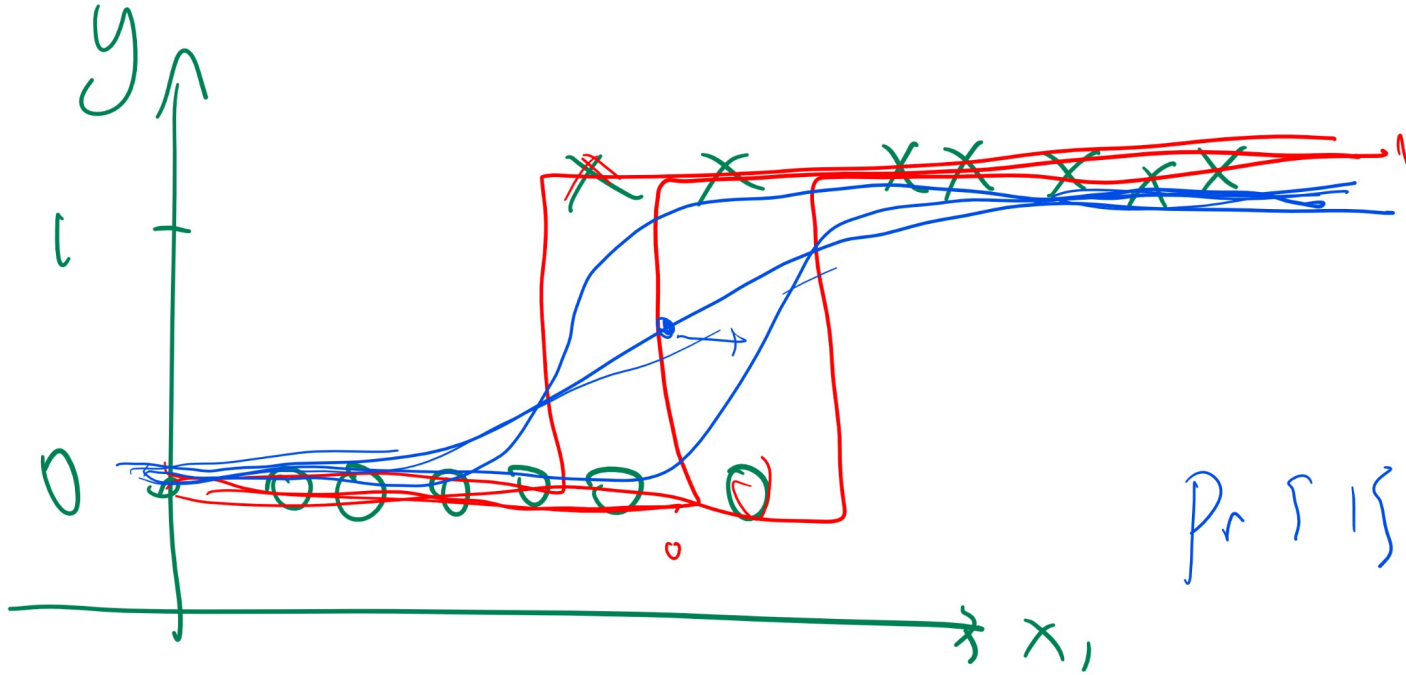
- When the activation function g is the sigmoid function, the problem is called logistic regression.

$$g(\square) = \frac{1}{1 + \exp(-\square)}$$

- This g is also called the logistic function.

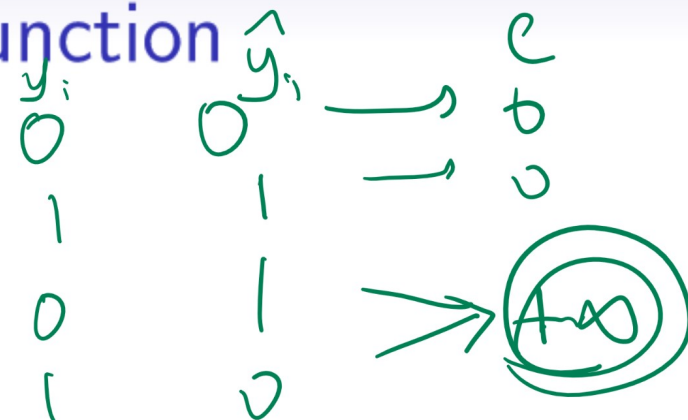
Sigmoid Function Diagram

Motivation



Cross Entropy Loss Function

Motivation



- The cost function used for logistic regression is usually the log cost function.

$$C(f) = - \sum_{i=1}^n (y_i \log(f(x_i)) + (1 - y_i) \log(1 - f(x_i)))$$

$y_i = 0$ $\hat{y}_i = 0$ $C = 0$
 $1 \cdot \log(0) = -\infty$

- It is also called the cross-entropy loss function.

Logistic Regression Objective

Motivation

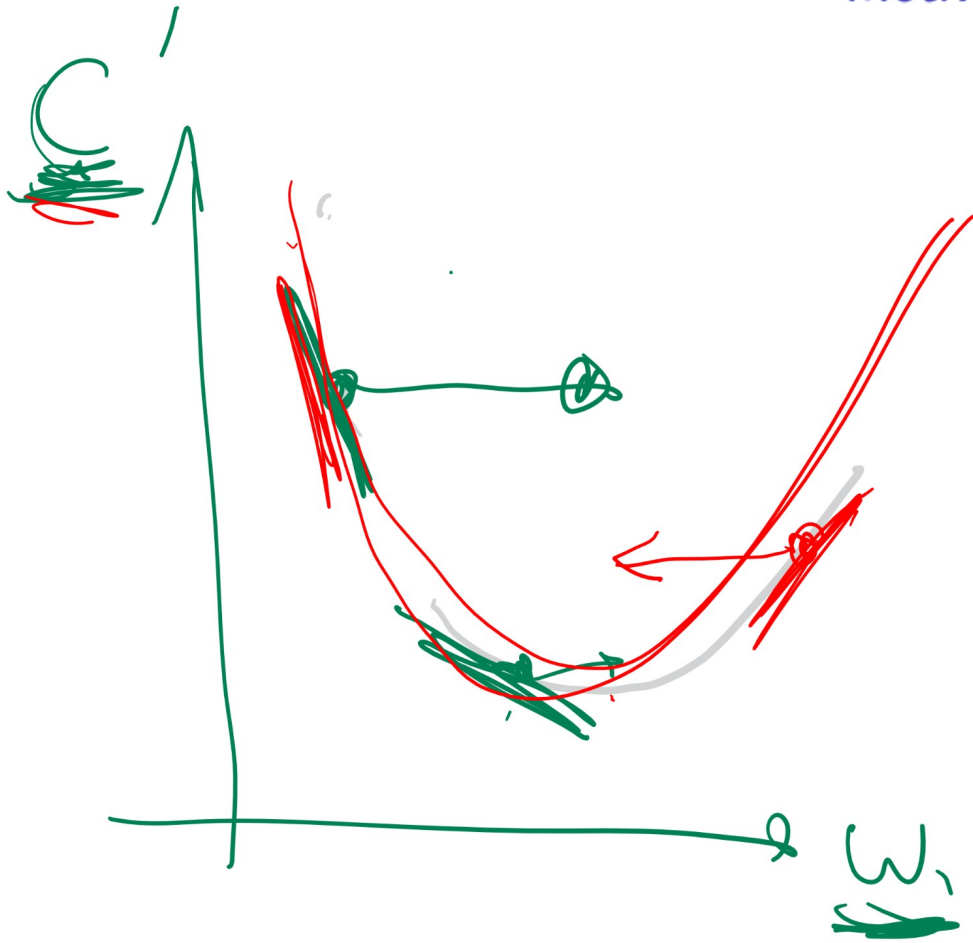
- The logistic regression problem can be summarized as the following.

$$(\hat{w}, \hat{b}) = \arg \min_{w, b} \sum_{i=1}^n (y_i \log(a_i) + (1 - y_i) \log(1 - a_i))$$

where $a_i = \frac{1}{1 + \exp(-z_i)}$ and $z_i = w^T x_i + b$

Optimization Diagram

Motivation



slope

large

neg

small

neg

pos

here
w ↑
by a lot

here
w ↑
by a little

decr
w ↓

Logistic Regression

Description

- Initialize random weights.
- Evaluate the activation function.
- Compute the gradient of the cost function with respect to each weight and bias.
- Update the weights and biases using gradient descent.
- Repeat until convergent.

Gradient Descent Intuition

Definition

Demo

- If a small increase in w_1 causes the distances from the points to the regression line to decrease: increase w_1 .
- If a small increase in w_1 causes the distances from the points to the regression line to increase: decrease w_1 .
- The change in distance due to change in w_1 is the derivative.
- The change in distance due to change in $\begin{bmatrix} w \\ b \end{bmatrix}$ is the gradient.

Gradient

Definition

- The gradient is the vector of derivatives.

- The gradient of

$$f(x_i) = w^T x_i + b = w_1 x_{i1} + w_2 x_{i2} + \dots + w_m x_{im} + b \text{ is:}$$

$$\nabla_w f = \begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \frac{\partial f}{\partial w_2} \\ \dots \\ \frac{\partial f}{\partial w_m} \end{bmatrix} = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \dots \\ x_{im} \end{bmatrix} = x_i$$

$$\nabla_b f = 1$$

Chain Rule

Definition

- The gradient of

$f(x_i) = g(w^T x_i + b) = g(w_1 x_{i1} + w_2 x_{i2} + \dots + w_m x_{im} + b)$
 can be found using the chain rule.

$$\nabla_w f = g'(w^T x_i + b) x_i$$

$$\nabla_b f = g'(w^T x_i + b) \mathbf{1}$$

Handwritten notes:
 $\frac{\partial}{\partial w} w^T x + b$
 x_i
 $\mathbf{1}$

- In particular, for the logistic function g :

$$g(\square) = \frac{1}{1 + \exp(-\square)}$$

$$g'(\square) = g(\square) (1 - g(\square))$$

- 8 - 4

Logistic Gradient Derivation 1

Definition

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

$$g'(z) = \frac{\cancel{1} \cdot \exp(-z) \cdot \cancel{1}}{(1 + \exp(-z))^2}$$

$$\frac{1}{x} \rightarrow -\frac{1}{x^2}$$

$$g'(z) =$$

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

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

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

$$\nabla_{\omega} f = a_i =$$

$$\left(\frac{g(z)}{a_i} \cdot (1 - \frac{g(z)}{a_i}) \right) \cdot X_i$$

Logistic Gradient Derivation 2

Definition

$$\nabla_w C = - \sum_{i=1}^n y_i \log a_i + (1 - y_i) \log(1 - a_i)$$

$\nabla_w a_i = a_i(1 - a_i) \cdot x_i$

$$\frac{\partial C}{\partial a_i} = y_i \frac{1}{a_i} - (1 - y_i) \frac{1}{1 - a_i}$$

$$= \frac{y_i(1 - a_i) - (1 - y_i)a_i}{a_i(1 - a_i)}$$

one term

$$\nabla_w C = \sum_{i=1}^n (a_i - y_i) x_i$$

Gradient Descent Step

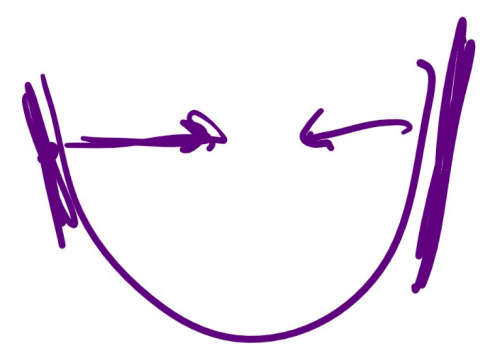
Definition

- For logistic regression, use chain rule twice.

$$w = w - \alpha \sum_{i=1}^n (a_i - y_i) x_i$$

$$b = b - \alpha \sum_{i=1}^n (a_i - y_i)$$

$$a_i = g(w^T x_i + b), \quad g(\square) = \frac{1}{1 + \exp(-\square)}$$



- α is the learning rate. It is the step size for each step of gradient descent.

Perceptron Algorithm

Definition

- Update weights using the following rule.

$$w = w - \alpha (a_i - y_i) x_i$$

$$b = b - \alpha (a_i - y_i)$$

$$a_i = \mathbb{1}_{\{w^T x_i + b \geq 0\}}$$

0, 1

Learning Rate Diagram

Definition

Logistic Regression, Part 1

Algorithm

- Inputs: instances: $\{x_i\}_{i=1}^n$ and $\{y_i\}_{i=1}^n$
- Outputs: weights and biases: w_1, w_2, \dots, w_m and b
- Initialize the weights.

$$w_1, \dots, w_m, b \sim \text{Unif} \left[\overset{\sim 1}{\cancel{0}}, 1 \right]$$

- Evaluate the activation function.

$$a_i = g(w^T x_i), g(\boxed{\cdot}) = \frac{1}{1 + \exp(-\boxed{\cdot})}$$

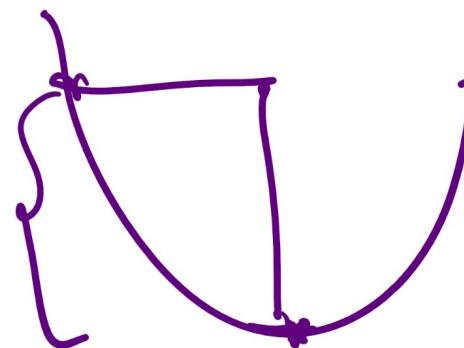
Logistic Regression, Part 2

Algorithm

- Update the weights and bias using gradient descent.

$$w = w - \alpha \sum_{i=1}^n (a_i - y_i) x_i$$

$$b = b - \alpha \sum_{i=1}^n (a_i - y_i)$$



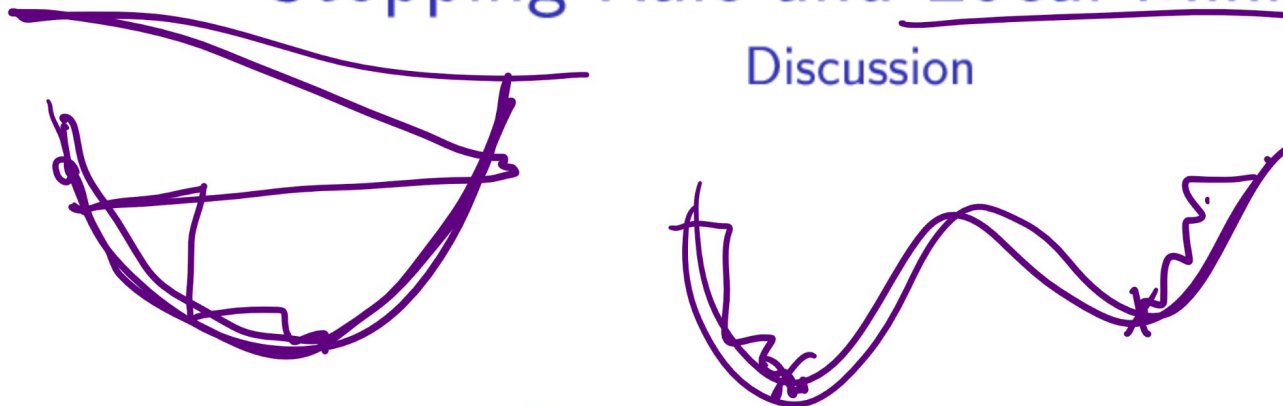
- Repeat the process until convergent.

$$|C - C^{\text{prev}}| < \epsilon$$

0.0001

Stopping Rule and Local Minimum

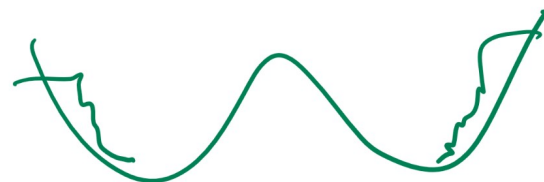
Discussion



- Start with multiple random weights.
- Use smaller or decreasing learning rates. One popular choice is $\frac{\alpha}{\sqrt{t}}$, where t is the iteration count.
- Use the solution with the lowest C.

Other Non-linear Activation Function

Discussion



- Activation function: $g(\cdot) = \tanh(\cdot) = \frac{e^{\cdot} - e^{-\cdot}}{e^{\cdot} + e^{-\cdot}}$

- Activation function: $g(\cdot) = \arctan(\cdot)$

- Activation function (rectified linear unit): $g(\cdot) = \underbrace{\cdot}_{\text{ReLU}} \mathbb{1}_{\{\cdot \geq 0\}}$

- All these functions lead to objective functions that are convex and differentiable (almost everywhere). Gradient descent can be used.

Convexity

Discussion

- If a function is convex, gradient descent with any initialization will converge to the global minimum.
- If a function is not convex, gradient descent with different initializations may converge to different local minima.
- A twice differentiable function is convex if and only if its second derivative is non-negative.
- In the multivariate case, it means the Hessian matrix is positive semidefinite.

$$f''(w) \geq 0$$

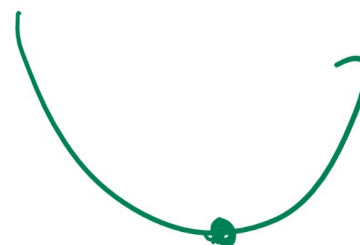
Positive Semidefinite

Discussion

$$f(x_1, x_2) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix}$$

- Hessian matrix is the matrix of second derivatives:

$$H : H_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$$



- A matrix H is positive semidefinite if $x^T H x \geq 0 \forall x \in \mathbb{R}^n$.
- A symmetric matrix is positive semidefinite if and only if all of its eigenvalues are non-negative.

$$H x = \lambda x$$

matrix eigen vector scalar eigen value

$$x^T H x = \lambda x^T x \geq 0$$

Convex Function Example 1

Discussion

$$f(x, y) = x^2 + xy + y^2$$

$$Hess = \begin{bmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial y \partial x} & \frac{\partial^2 f}{\partial y^2} \end{bmatrix} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$$

$$\nabla f = \begin{pmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{pmatrix} = \begin{pmatrix} 2x + y \\ 2y + x \end{pmatrix}$$

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 3 \\ 3 \end{bmatrix} \approx \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Convex Function Example 2

Discussion

$$\begin{aligned} \underbrace{\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}}_{\lambda_1} \begin{bmatrix} 1 \\ -1 \end{bmatrix} &= \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \underbrace{0}_{\lambda_2} \begin{bmatrix} 1 \\ -1 \end{bmatrix} \\ \underbrace{\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}}_{\lambda_1} \begin{bmatrix} -1 \\ 1 \end{bmatrix} &= \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \underbrace{0}_{\lambda_2} \begin{bmatrix} -1 \\ 1 \end{bmatrix} \end{aligned}$$

$$\lambda_1, \lambda_2 = 1, 3 > 0$$

pos. semi. det

Convex Functions

Quiz

Definiteness

Quiz