

CS5242 Neural Networks and Deep Learning

Lecture 02: Shallow Networks

Wei WANG

cs5242@comp.nus.edu.sg

change log:

- slide 32, loss function of the vectorized multi-label multi-classification model
- removed the overfitting part



NUS
National University
of Singapore

School of
Computing

Agenda

- Recap of last lecture
- Classification
 - Logistic regression
 - binary cross-entropy
 - Multinomial regression (Softmax regression)
- Overfitting and underfitting

Recap

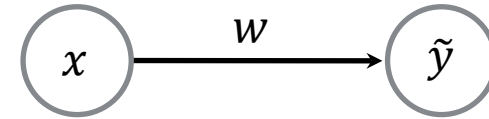
- Linear regression
 - Univariate: single feature $x \in R$
 - Multivariate: multiple features $\mathbf{x} \in R^m$
 - Linear transformation: $\tilde{y} = \mathbf{w}^T \mathbf{x}$
 - Loss: measure the difference between the prediction and ground truth
 - Training is to optimize (i.e., minimize) the loss w.r.t parameters (\mathbf{w})
- Gradient descent algorithm
 - Minimize the target loss iteratively; for each iteration,
 - Compute the gradient of the average loss (over all training examples) w.r.t \mathbf{w}
 - Update \mathbf{w} in the **opposite** of the gradient direction

$$\mathbf{w} = \mathbf{w} - \alpha \frac{\partial J}{\partial \mathbf{w}}$$

The simplest case

- Univariate linear regression without bias

- $\tilde{y} = wx$



- Squared error as the loss function

- $L(x, y|w) = \frac{1}{2} ||\tilde{y} - y||^2 = \frac{1}{2} (\tilde{y} - y)^2 = \frac{1}{2} (wx - y)^2$

- Training objective

- tune w to minimize the average loss over the training dataset

- $J(w) = \frac{1}{m} \sum_{i=1}^m L(x^{(i)}, y^{(i)}|w) = \frac{1}{2m} \sum_{i=1}^m (wx^{(i)} - y^{(i)})^2$

Gradient descent algorithm

Randomly initialize w

Repeat

$$J = 0$$

For each sample $x^{(i)}, y^{(i)}$ ($i=1, 2, \dots, m$)

a. Forward pass $\tilde{y}^{(i)} = wx^{(i)}$

b. Accumulate Loss $J += \frac{1}{m} L^{(i)} = \frac{1}{2m} \|\tilde{y}^{(i)} - y^{(i)}\|^2 = \frac{1}{2m} (wx^{(i)} - y^{(i)})^2$

c. Accumulate gradient $\frac{\partial J}{\partial w} += \frac{1}{m} \frac{\partial L^{(i)}}{\partial w} = \frac{1}{2m} (wx^{(i)} - y^{(i)})x^{(i)}$

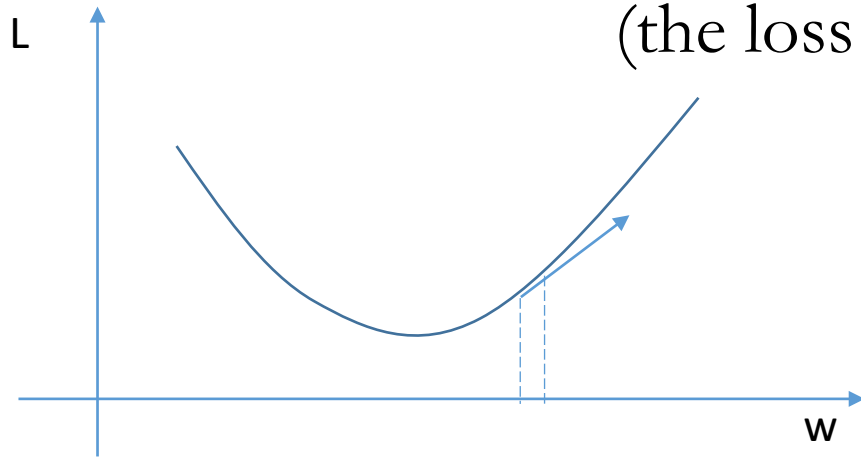
Update $w = w - \alpha \frac{\partial J}{\partial w}$ *move in the opposite direction of gradient*

Understanding Derivatives

$\frac{d\mathcal{L}}{dw}$...is the rate at which **this** will change...

$$\mathcal{L} = \frac{1}{2}(y - \hat{y})^2$$

(the loss function)



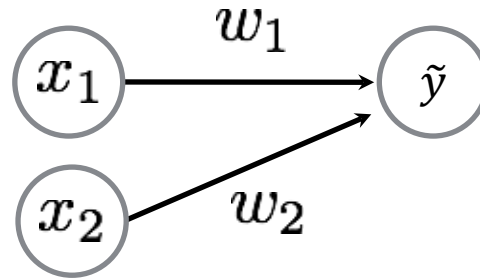
... per unit change of **this**

$$y = wx$$

(the weight parameter)

The second simplest case

- Multivariate linear regression with two features
 - the bias term is skipped
 - $\tilde{y} = \mathbf{w}^T \mathbf{x} = w_1 x_1 + w_2 x_2$



Gradient descent algorithm

Random initialize w_1, w_2

Repeat

$J = 0$

For each sample $\mathbf{x}^{(i)}, y^{(i)}$

a. Forward pass

b. Accumulate Loss

c. Accumulate partial gradient (derivative) $\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2}$ *what are ?*

Update

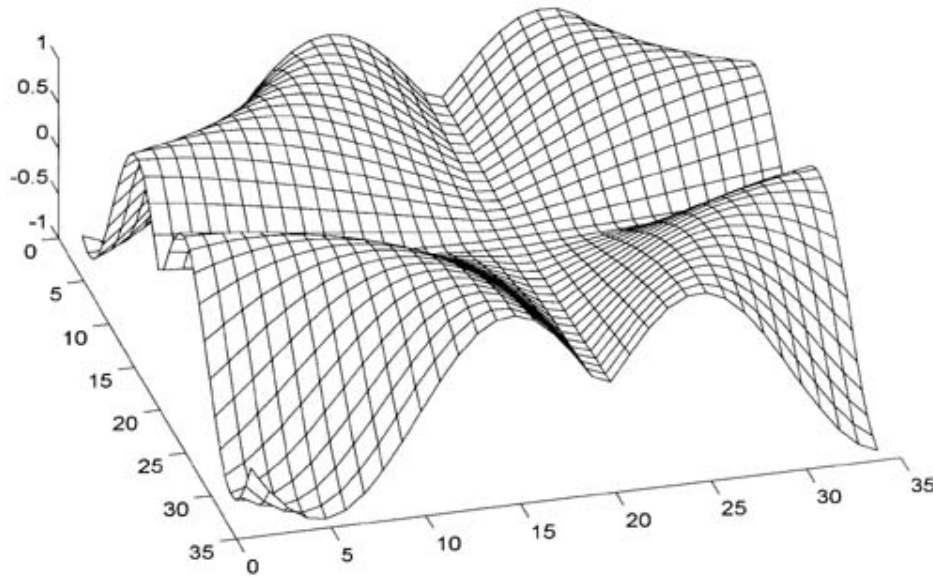
$$w_1 = w_1 - \alpha \frac{\partial J}{\partial w_1}$$

$$w_2 = w_2 - \alpha \frac{\partial J}{\partial w_2}$$

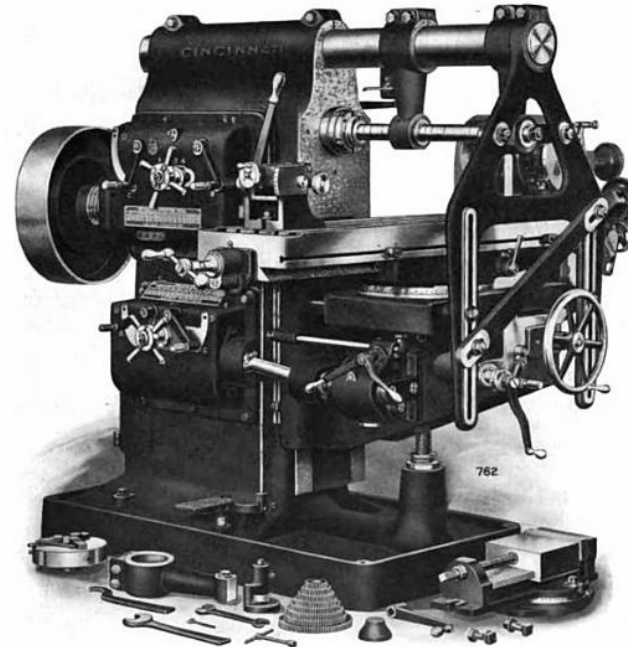
move in opposite direction of partial derivatives

(Partial) Derivatives

Two ways to think about them:

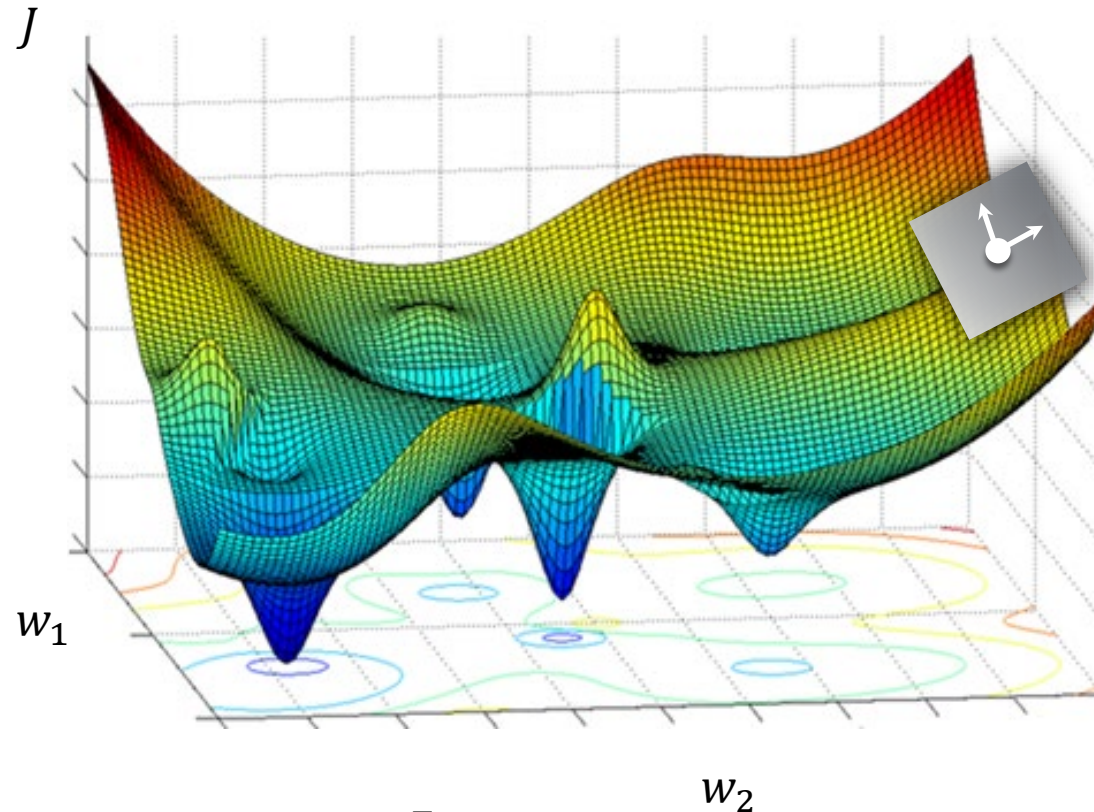


Slope of a function



Knobs on a machine

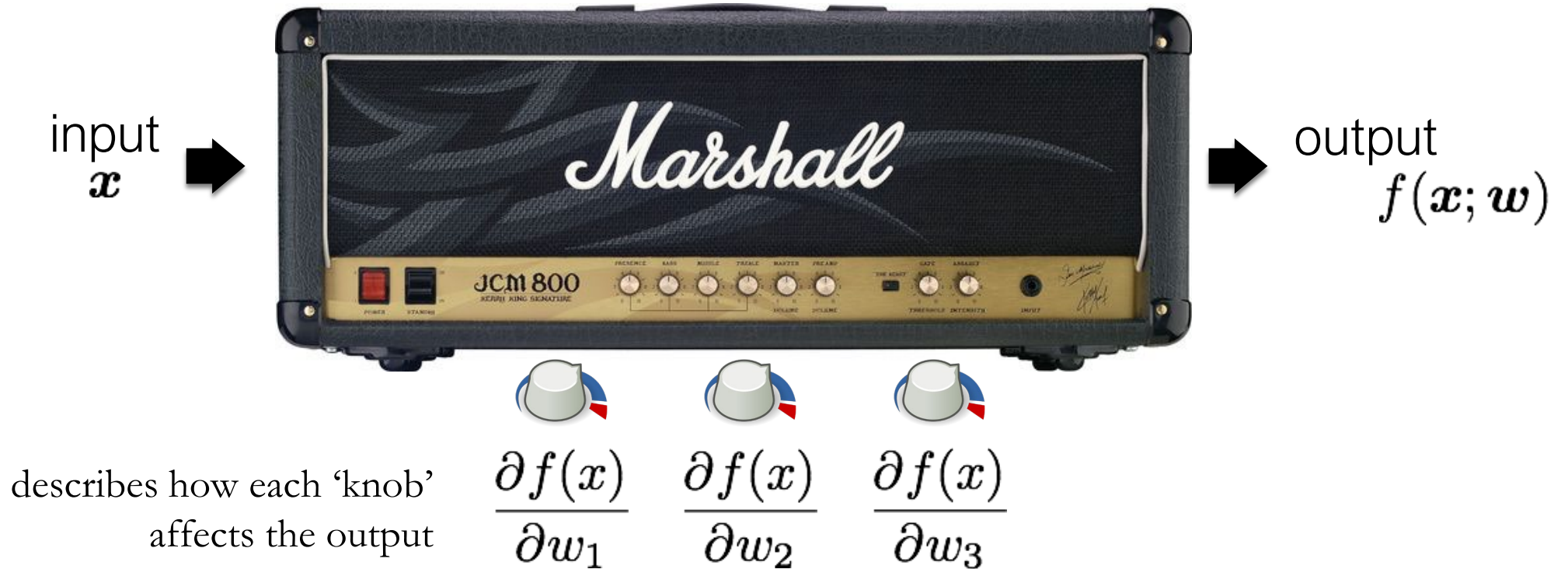
Slope of a function



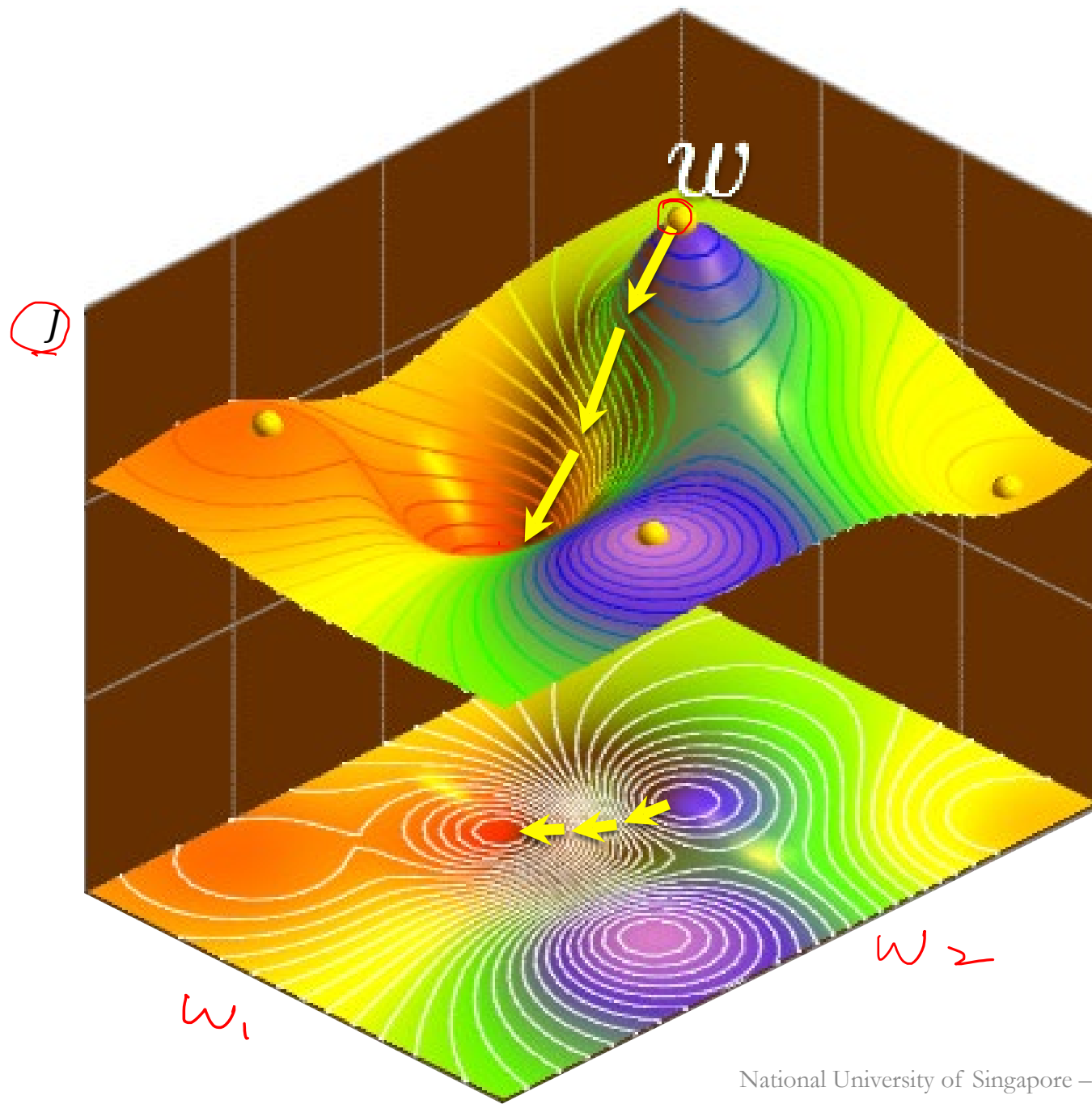
$$\frac{\partial J}{\partial \mathbf{w}} = \left[\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2} \right]^T$$

describes the slope around a point

Knobs on a machine



small change in parameter Δw_1 → output will change by $\frac{\partial f(x)}{\partial w_1} \Delta w_1$



Gradient Descent:
given a fixed-point on a
function, move in the
direction opposite of the
gradient.

$$w_i = w_i - \alpha \frac{\partial J}{\partial w_i}$$

Recap

- Vectorization & denominator layout
- Let $\Delta/\blacksquare/\nabla$ be a scalar, vector or matrix.
- $\frac{\partial \Delta}{\partial \blacksquare}$
 - If Δ is scalar, then the result shape is the same as \blacksquare
 - If \blacksquare is scalar, then the result shape is the same as Δ^T
 - If both are vectors, then the result is a matrix with rows decided by \blacksquare and columns by Δ
- $\Delta = g(\blacksquare), \blacksquare = u(\nabla)$
 - $\frac{\partial \Delta}{\partial \nabla}$ is the multiplication between $g'(\blacksquare)$ and $u'(\nabla)$, but
 - Need to arrange the order and transpose to make sure the result's shape matches $\frac{\partial \Delta}{\partial \nabla}$

Gradient table

(denominator layout)

Shape of $\frac{\partial y}{\partial x}$	Scalar y	Vector \mathbf{y} (m, 1)	Matrix \mathbf{Y} (m, n)
Scalar x	1	(1, m)	(n, m)
Vector \mathbf{x} (n, 1)	(n, 1)	(n, m)	
Matrix \mathbf{X} (p, q)	(p, q)		

Shape check:

For denominator layout, the shape of the gradient of a scalar w.r.t a variable v is the same as the shape of v , where v could be a scalar, vector, or matrix

Gradient table: vector by vector (denominator layout)

$y =$	a	x	Ax	$x^T A$
$\frac{\partial y}{\partial x} =$	0	I	A^T	A

Gradient table: vector by vector (denominator layout)

$y =$	au $u = u(x)$	vu $v = v(x), u = u(x)$	$v + u$ $v = v(x), u = u(x)$	Au $u = u(x)$	$g(u)$ $u = u(x)$
$\frac{\partial y}{\partial x} =$	$a \frac{\partial u}{\partial x}$	$v \frac{\partial u}{\partial x} + \frac{\partial v}{\partial x} u^T$	$\frac{\partial v}{\partial x} + \frac{\partial u}{\partial x}$	$\frac{\partial u}{\partial x} A^T$	$\frac{\partial u}{\partial x} \frac{\partial g(u)}{\partial u}$

Gradient table: scalar by vector

(denominator layout)

$y =$	a	$u^T v$ $v = v(x), u = u(x)$	$g(u)$ $u = u(x)$	$x^T A x$
$\frac{\partial y}{\partial x} =$	0	$\frac{\partial v}{\partial x} u + \frac{\partial u}{\partial x} v$	$\frac{\partial g(u)}{\partial u} \frac{\partial u}{\partial x}$	$(A + A^T)x$

Vectorization

- Consider a single instance

$$J(\mathbf{w}) = L(\mathbf{x}, y | \mathbf{w}) = \frac{1}{2} \|\tilde{y} - y\|^2 = \frac{1}{2} (\mathbf{w}^T \mathbf{x} - y)^2$$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = ?$$

$$J(\mathbf{w}) = \frac{1}{2} z^2, \text{ where } z = \mathbf{w}^T \mathbf{x} - y$$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \frac{\partial J}{\partial z} \frac{\partial z}{\partial \mathbf{w}} = z \mathbf{x} = (\mathbf{w}^T \mathbf{x} - y) \mathbf{x}$$

Shape Check!

Vectorization

- Consider multiple examples $\{(\mathbf{x}^{(i)}, y^{(i)})\}, i=1, 2, \dots, m$

$$X = \begin{pmatrix} \mathbf{x}^{(1)T} \\ \mathbf{x}^{(2)T} \\ \vdots \\ \mathbf{x}^{(m)T} \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{pmatrix} \quad \tilde{\mathbf{y}} = X\mathbf{w}$$

$$J(\mathbf{w}) = \frac{1}{2m} \|\tilde{\mathbf{y}} - \mathbf{y}\|^2 = \frac{1}{2m} (\tilde{\mathbf{y}} - \mathbf{y})^T (\tilde{\mathbf{y}} - \mathbf{y}) = \frac{1}{2m} \mathbf{u}^T \mathbf{u} \quad (\mathbf{u} = \tilde{\mathbf{y}} - \mathbf{y})$$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{2m} \left(\frac{\partial \mathbf{u}}{\partial \mathbf{w}} \mathbf{u} + \frac{\partial \mathbf{u}}{\partial \mathbf{w}} \mathbf{u} \right) = \frac{1}{m} \mathbf{X}^T \mathbf{u} = \frac{1}{m} \mathbf{X}^T (\tilde{\mathbf{y}} - \mathbf{y})$$

Shape Check!

Classification

Regression vs. classification, cross-entropy loss

Multi-class, multi-label classification

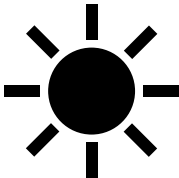
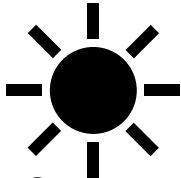


Multi-class, single-label classification

Regression VS Classification

Q: What is the difference?

Quantity vs. Label:
regression maps to a continuous domain, while classification maps to a finite set.

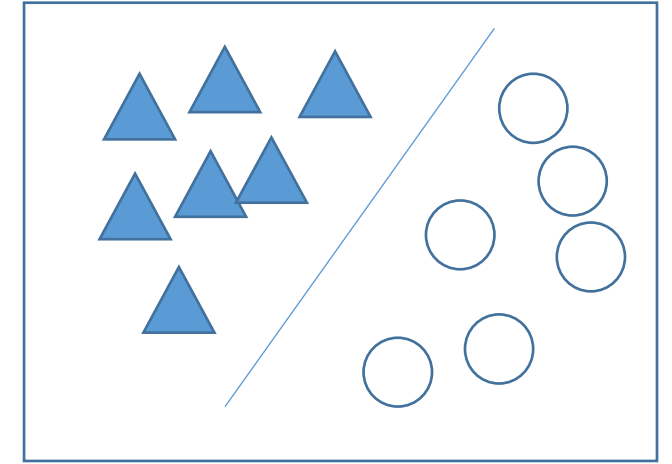
- Regression: What's the temperature of tomorrow?
- Classification (Binary): Is it sunny tomorrow?
- Classification (Multi-class): Is it sunny, cloudy, or rainy?

		  		
		Sunny	Rainy	Cloudy
today	1	Monday	1	0
tomorrow	0	Tuesday	0	1
		Wednesday	0	0

Q: How many mm of rain makes it a "rainy" day?
How many hours of sunshine makes it a "sunny" day?

Often, classification labels must be derived from continuous values (measured or regressed).

From regression to classification



- Thresholding (Perceptron)
- $\tilde{y} = \begin{cases} 1, & \text{if } \mathbf{w}^T \mathbf{x} > c \\ 0, & \text{else} \end{cases}$
- How to set the threshold c ?

Learn it as a part of learning weights.

$$\mathbf{w}^T \mathbf{x} > c$$

$$x_1 w_1 + x_2 w_2 \dots + x_m w_m + b > c$$

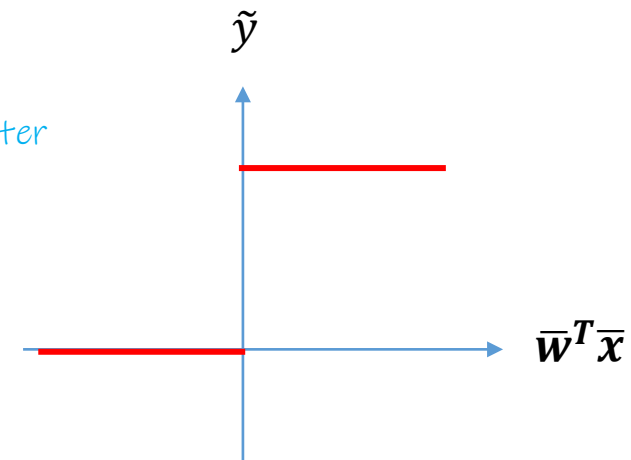
$$x_1 w_1 + x_2 w_2 \dots + x_m w_m + \underbrace{(b - c)}_{\text{Merge } c \text{ as part of the offset / } b \text{ parameter}} > 0$$

Merge c as part of the offset / b parameter

$$\bar{\mathbf{w}}^T \bar{\mathbf{x}} > 0$$

$$\tilde{y} = \begin{cases} 1, & \text{if } \bar{\mathbf{w}}^T \bar{\mathbf{x}} > 0 \\ 0, & \text{else} \end{cases}$$

What is $\bar{\mathbf{w}}$, $\bar{\mathbf{x}}$?



Logistic regression

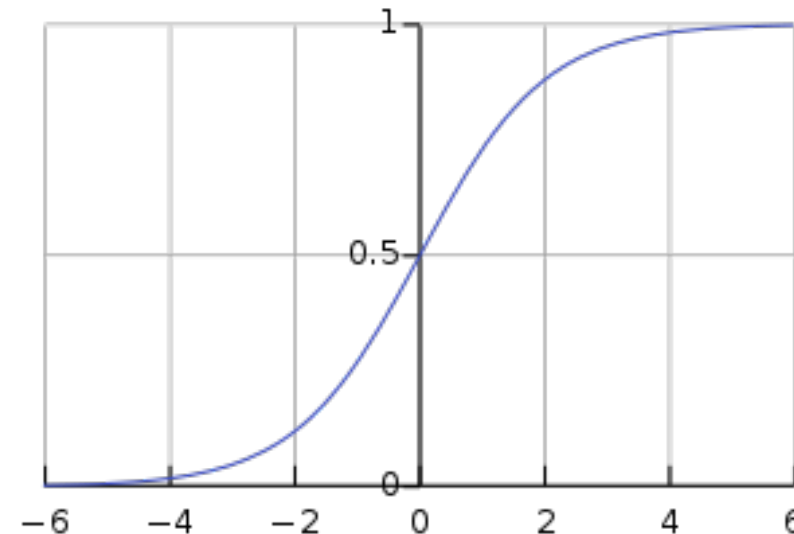
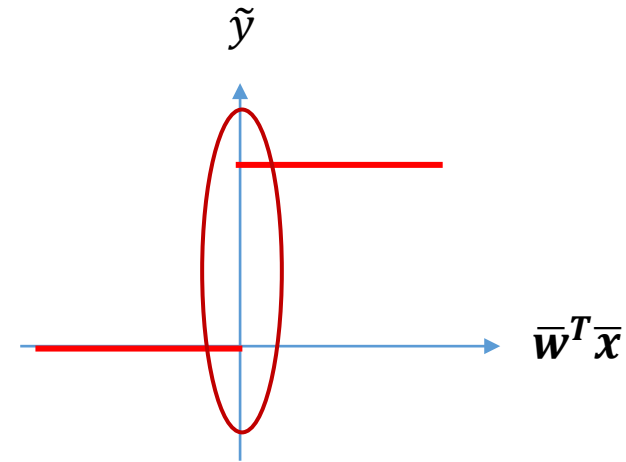
$$\tilde{y} = \begin{cases} 1, & \text{if } \mathbf{w}^T \mathbf{x} > 0 \\ 0, & \text{else} \end{cases}$$

Gradient is always 0, so cannot learn anything!

Logistic function: $p = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}}}$

- Range is within $[0, 1]$
- Possible interpretation: probability of the label being 1
- Logistic function sometimes also referred to as a sigmoid function

Only piecewise differentiable



This Photo by Unknown Author is licensed under CC BY-SA

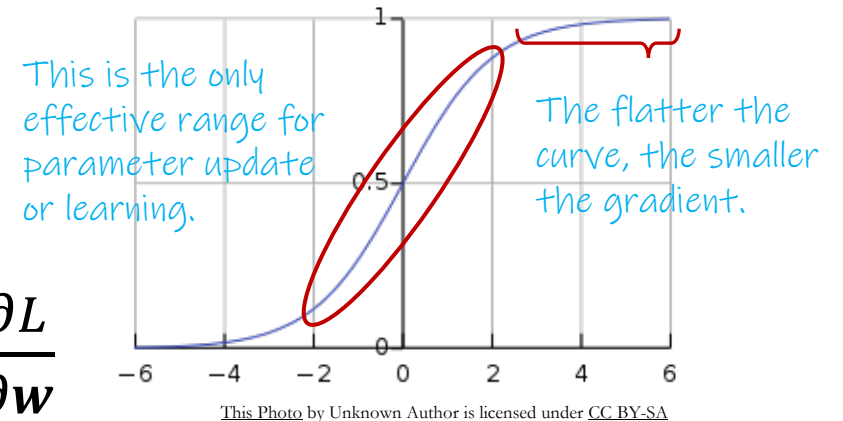
Gradient Vanishing

How should we learn the weights of the logistic function? Start with a simple L2 loss.

- $L(\mathbf{x}, y) = \frac{1}{2} ||\sigma(\mathbf{w}^T \mathbf{x}) - y||^2, \frac{\partial L}{\partial \mathbf{w}}?$
 - denote $z = \mathbf{w}^T \mathbf{x}, L = \frac{1}{2} (\sigma(z) - y)^2$

HWQ: derive the following gradient result.

- $\frac{\partial L}{\partial \mathbf{w}} = (\sigma(z) - y) * \sigma(z)(1 - \sigma(z)) \mathbf{x}$
- If $\sigma \approx 0$ or $1, \frac{\partial \sigma}{\partial z} \approx 0 \rightarrow \frac{\partial L}{\partial \mathbf{w}} \approx \mathbf{0}$ gradient vanishing
- Impact: training gets stuck since $\mathbf{w} = \mathbf{w} - \alpha * \frac{\partial L}{\partial \mathbf{w}}$



Cross-entropy loss

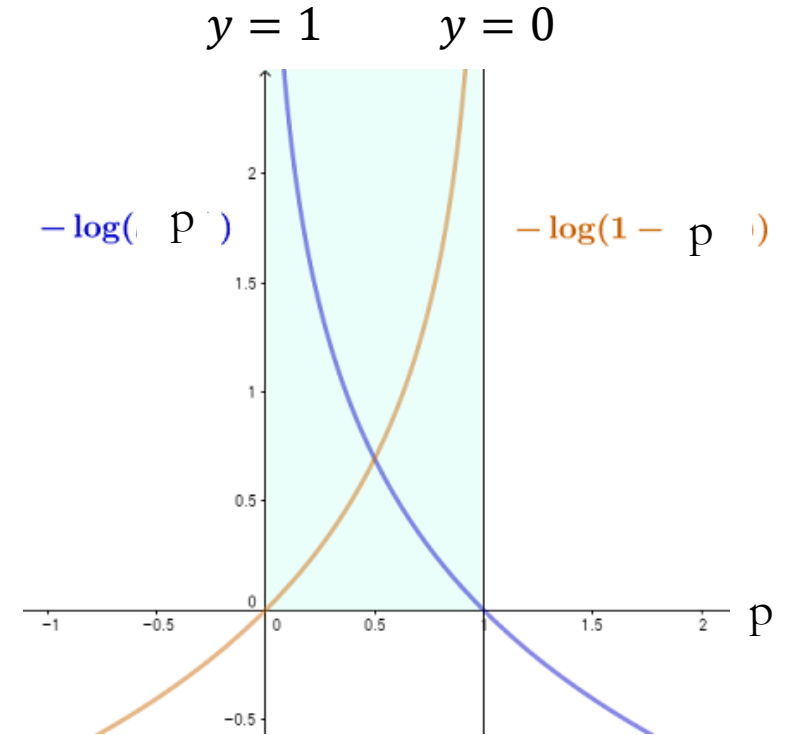
- Denote $p = \sigma(z), z = \mathbf{w}^T \mathbf{x}$
- Instead of using an L2 loss, we use the following:
- $L_{ce}(x, y) = \underbrace{-y \log p}_{\text{This term goes to 0 if ground truth label is 0}} - \underbrace{(1 - y) \log(1 - p)}_{\text{This term goes to 0 if ground truth label is 1}}$

This term goes to 0 if
ground truth label is 0

This term goes to 0 if
ground truth label is 1

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

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



This Photo by Unknown Author is licensed under CC BY-SA

What is the intuition behind this loss? Does it actually help us learn the right weights?

Cross-entropy loss explanation

Consider the probability of a classifier being correct.

$$P(\text{correct}|\mathbf{x}) = \begin{cases} P(\tilde{y} = 1|\mathbf{x}), & \text{if } y = 1 \\ P(\tilde{y} = 0|\mathbf{x}), & \text{if } y = 0 \end{cases} \quad (\text{depends on the ground truth label } y)$$

$$= P(\tilde{y} = 1|\mathbf{x})^y P(\tilde{y} = 0|\mathbf{x})^{1-y} \quad \text{collapse cases into a single function}$$

Log-likelihood of our classifier being correct:

$$\log P(\text{correct}|\mathbf{x}) = y \log P(\tilde{y} = 1|\mathbf{x}) + (1 - y) \log P(\tilde{y} = 0|\mathbf{x})$$

We want to maximize this, i.e. to maximize the probability of our classifier being correct!

Objective equivalent to minimizing the negative log-likelihood

$$\min -\log P(\text{correct}|\mathbf{x}) = \min -y \log P(\tilde{y} = 1|\mathbf{x}) - (1 - y) \log P(\tilde{y} = 0|\mathbf{x})$$

Note that so far, this is general and that we have not made any assumptions about the classifier itself, i.e. the specific form of $P(\tilde{y}|\mathbf{x})$

Cross-entropy loss explanation

$$P(\tilde{y} = 1|\mathbf{x}) = p = \sigma(z)$$

Because range of logistic is between 0-1, we adopt p as the probability of the of x having a label.

$$P(\tilde{y} = 0|\mathbf{x}) = 1 - p$$

Problem is binary, equate probability of having label 0 as the complement.

Minimizing negative log likelihood:

$$\min -\log P(\text{correct}|\mathbf{x}) = \min -y \log P(\tilde{y} = 1|\mathbf{x}) - (1 - y) \log P(\tilde{y} = 0|\mathbf{x})$$

$$= \min \underbrace{-y \log p - (1 - y) \log (1 - p)}_{L_{ce}(\mathbf{x}, y)}$$

Substituting the logistic function for $P(\tilde{y}|\mathbf{x})$

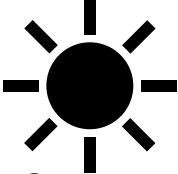


The cross-entropy loss minimizes the classifier's log-likelihood.

Cross-entropy loss gradient

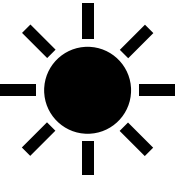


$$\begin{aligned}\frac{\partial L_{ce}}{\partial \mathbf{w}} &= \frac{\partial L_{ce}}{\partial z} \frac{\partial z}{\partial \mathbf{w}} = \frac{\partial L_{ce}}{\partial p} \frac{\partial p}{\partial z} \frac{\partial z}{\partial \mathbf{w}} \\ &= \left(-\frac{y}{p} + \frac{1-y}{1-p} \right) * p * (1-p) \mathbf{x} \\ &= (p - y) \mathbf{x}\end{aligned}$$

Multi-class classification

- Single-label

			
	Sunny	Rainy	Cloudy
Monday	1	0	0
Tuesday	0	1	0
Wednesday	0	0	1

- Multi-label: can belong to more than one class

			
Monday	1	0	1
Tuesday	0	1	1

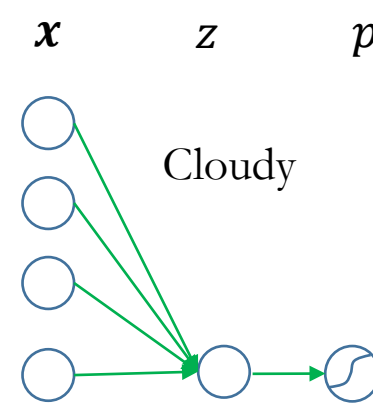
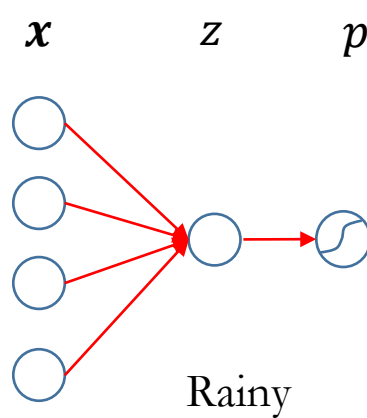
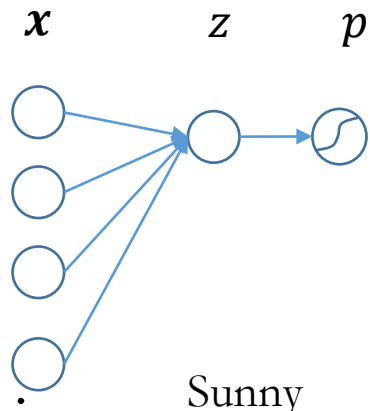
Sunny & Cloudy
Rainy & Cloudy

Applications

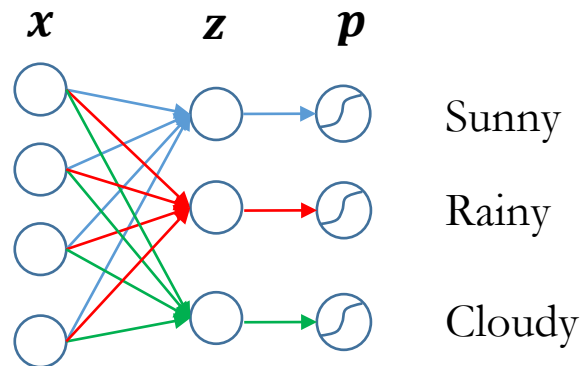
- Multi-class image classification
 - Classify each image into one of the class
 - [MNIST dataset](#)
 - $\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$
 - What is \mathbf{x} ? i.e., how to represent an image
 - [Cifar10 dataset](#)
 - {Dog, Cat, Horse, Ship, Truck, Frog, Deer, Bird, Automobile, Airplane}
- Multi-class document classification
 - 20 Newsgroups, {hardware, autos, space, etc}

Multi-class multi-label classification

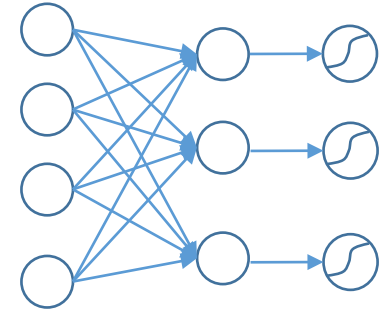
- Binary classification for each label



- All in one network



Multi-class multi-label classification

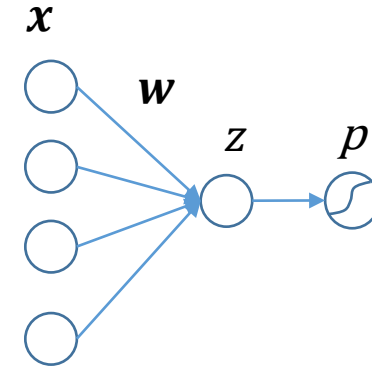


- For binary classification
 - $z = \mathbf{w}^T \mathbf{x} + b, \quad p = \sigma(z), \quad L_{ce} = -y \log p - (1 - y) \log(1 - p)$
 - $\mathbf{x} \in R^n, \mathbf{w} \in R^n, b \in R, p \in R$
- For multi-class, the i^{th} class
 - We assume that the classes are independently conditioned on the input
 - $z_i = \mathbf{W}_i \mathbf{x} + b_i, \quad p_i = \sigma(z_i), \quad L_{ce} = -y_i \log p_i - (1 - y_i) \log(1 - p_i)$
 - $\mathbf{x} \in R^n, \mathbf{W}_i \in R^{1 \times n}, b_i \in R^1, p_i \in R^1$
- For multi-class, multi-label (vectorized form)
 - $\mathbf{z} = \mathbf{W} \mathbf{x} + \mathbf{b}, \quad \mathbf{p} = \sigma(\mathbf{z}), \quad L_{ce} = -\mathbf{y}^T \log \mathbf{p} - (\mathbf{1} - \mathbf{y})^T \log(\mathbf{1} - \mathbf{p})$
 - $\mathbf{x} \in R^n, \mathbf{W} \in R^{k \times n}, \mathbf{b} \in R^k, \mathbf{p} \in R^k, \mathbf{y} \in \{0, 1\}^k$

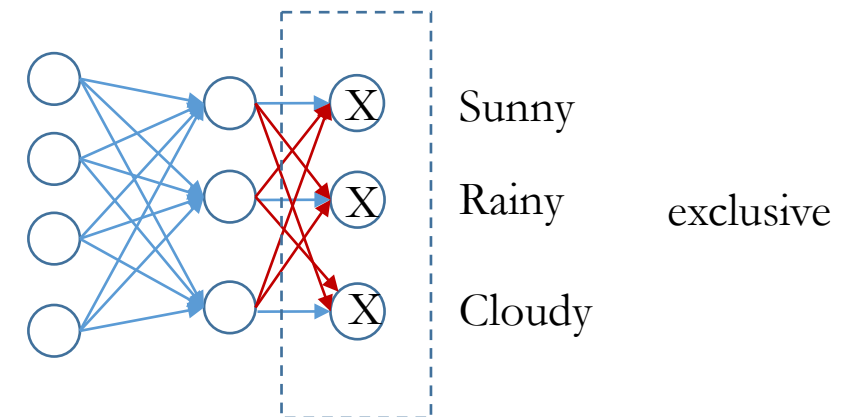
Probability of belonging to class i is independent of belonging to class j , once conditioned on the input evidence

Multi-class single-label classification

- Choose one label from multiple classes
 - Exclusive
 - If $p(\text{sunny})$ is large, then $p(\text{rainy}) + p(\text{cloudy})$ small
- How to enforce this constraint?
 - $p(\text{sunny}) + p(\text{rainy}) + p(\text{cloudy}) = 1$
 - $\sum_{i=1} p_i = 1$



Our previous model had connections only from each z_i to p_i ; now we add the red arrows to connect all z_i to p_i



Multi-class single-label classification

- Softmax regression or multinomial logistic regression

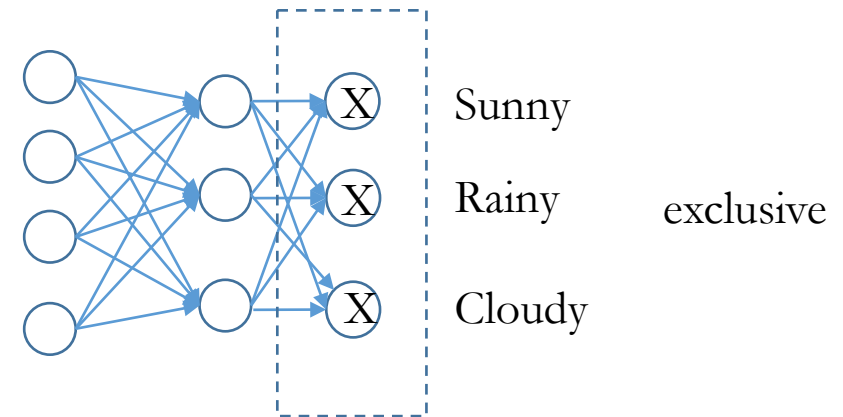
- $\mathbf{z} = \mathbf{W}\mathbf{x}, \mathbf{W} \in R^{k \times n},$

- $p_i = \frac{e^{z_i}}{\sum_j e^{z_j}} = \text{softmax}(z_i)$

- Then we have $\sum_i p_i = 1$

- z_i is called a logit

- $t = \underset{i}{\operatorname{argmax}} p_i; \tilde{y}_i = 1 \text{ if } i = t; \text{ else } 0;$



Multi-class single-label classification

- Loss function: $L_{ce}(\mathbf{x}, \mathbf{y}) = \sum_i -y_i \log p_i = -\mathbf{y}^T \log \mathbf{p}$
 - Each row of \mathbf{X} is the feature vector \mathbf{x}
 - Each row of \mathbf{Y} is the target one-hot vector \mathbf{y}

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}^{(1)T} \\ \mathbf{x}^{(2)T} \\ \dots \\ \mathbf{x}^{(N)T} \end{pmatrix} \qquad \mathbf{Y} = \begin{pmatrix} \mathbf{y}^{(1)T} \\ \mathbf{y}^{(2)T} \\ \dots \\ \mathbf{y}^{(N)T} \end{pmatrix}$$

$$\bullet \frac{\partial L_{ce}}{\partial \mathbf{W}} = ? \quad \frac{\partial L_{ce}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}}$$

HWQ: derive the gradient of L_{ce} w.r.t to \mathbf{W} .

Cross-entropy

- “Average number of bits needed to identify an event drawn from the set, if a coding scheme of the set is optimized for an estimated probability distribution p , rather than the “true” distribution q ”---
Wikipedia
- $H(q, p) = H(q) + D_{kl}(q||p) = -\sum_k q_k \log p_k$
 - $H(q)$ is a constant, i.e., the entropy of the true distribution
 - $D_{kl}(q||p)$ measures the difference between two distributions
 - Therefore, $H(q, p)$ measures the difference between two distributions
- For multi-class classification
 - \mathbf{p} is the artificial distribution to be optimized
 - \mathbf{y} is the true distribution (q), i.e., the label distribution

Summary

- Regression vs. classification as basic machine learning tasks
 - Using logistic regression to approximate the decision for classification
 - Logistic functions should be learned with cross-entropy and not L2 loss due to gradient vanishing problem
 - cross-entropy loss tries to minimize the difference between the output distribution and the (ground truth) label distribution

Homework Questions

1. Derive the gradient of the L2 loss wrt \mathbf{w} using a logistic regression function (slide 23)
2. Derive the gradient of the cross entropy loss (slide 27)
3. Derive the gradient of Lce w.r.t to \mathbf{W} (slide 34)

Next Lecture

Overfitting & Multi-layer perceptron & Back-propagation