# Linear Models: Perceptron, Logistic Regression

## CMSC 470

Marine Carpuat

# Linear Models for Multiclass Classification

$$\hat{y} = \arg\max_y \boldsymbol{\theta}^\mathsf{T} \mathbf{f}(\mathbf{x}, y)$$

Feature function representation

Weights

# Multiclass perceptron

---
**Algorithm 3** Perceptron learning algorithm

---
1: **procedure** $\text{PERCEPTRON}(\boldsymbol{x}^{(1:N)}, y^{(1:N)})$
2: $\quad\quad t \leftarrow 0$
3: $\quad\quad \boldsymbol{\theta}^{(0)} \leftarrow \mathbf{0}$
4: $\quad$ **repeat**
5: $\quad\quad\quad t \leftarrow t + 1$
6: $\quad\quad\quad$ Select an instance $i$
7: $\quad\quad\quad \hat{y} \leftarrow \text{argmax}_y \, \boldsymbol{\theta}^{(t-1)} \cdot \boldsymbol{f}(\boldsymbol{x}^{(i)}, y)$
8: $\quad\quad\quad$ **if** $\hat{y} \neq y^{(i)}$ **then**
9: $\quad\quad\quad\quad \boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)} + \boldsymbol{f}(\boldsymbol{x}^{(i)}, y^{(i)}) - \boldsymbol{f}(\boldsymbol{x}^{(i)}, \hat{y})$
10: $\quad\quad\quad$ **else**
11: $\quad\quad\quad\quad \boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)}$
12: $\quad$ **until** tired
13: $\quad$ **return** $\boldsymbol{\theta}^{(t)}$

---

# Properties of Linear Models we've seen so far

**Naïve Bayes**

- Batch learning

- Generative model p(x,y)

- Grounded in probability

- Assumes features are independent given class

- Learning = find parameters that maximize likelihood of training data

**Perceptron**

- Online learning

- Discriminative model score(y|x)

- Guaranteed to converge if data is linearly separable

- But might overfit the training set

- Error-driven learning

# Averaged Perceptron improves generalization
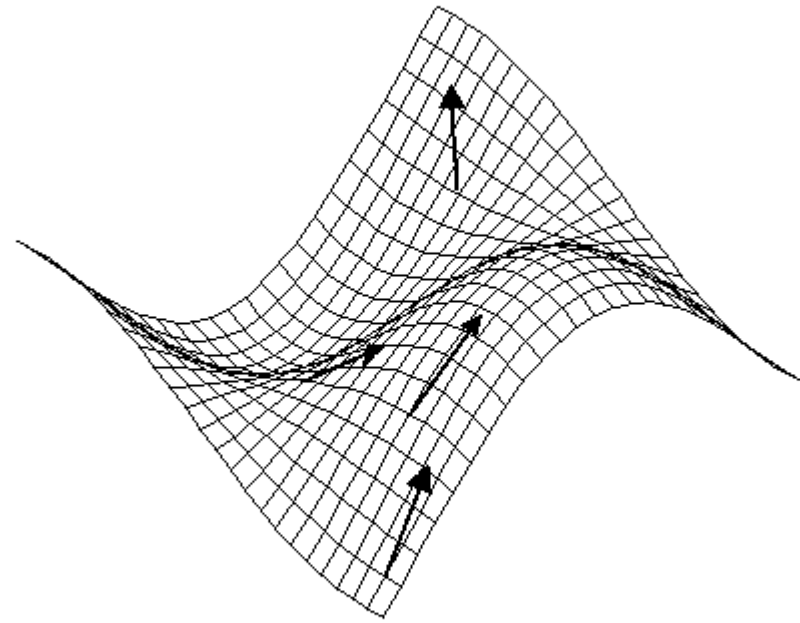
**Algorithm 4** Averaged perceptron learning algorithm

1: **procedure** AVG-PERCEPTRON($\boldsymbol{x}^{(1:N)}, \boldsymbol{y}^{(1:N)}$)
2: $\quad t \leftarrow 0$
3: $\quad \boldsymbol{\theta}^{(0)} \leftarrow 0$
4: $\quad$ **repeat**
5: $\quad\quad t \leftarrow t + 1$
6: $\quad\quad$ Select an instance $i$
7: $\quad\quad \hat{y} \leftarrow \text{argmax}_y \, \boldsymbol{\theta}^{(t-1)} \cdot \boldsymbol{f}(\boldsymbol{x}^{(i)}, y)$
8: $\quad\quad$ **if** $\hat{y} \neq y^{(i)}$ **then**
9: $\quad\quad\quad \boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)} + \boldsymbol{f}(\boldsymbol{x}^{(i)}, y^{(i)}) - \boldsymbol{f}(\boldsymbol{x}^{(i)}, \hat{y})$
10: $\quad\quad$ **else**
11: $\quad\quad\quad \boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)}$
12: $\quad\quad \boldsymbol{m} \leftarrow \boldsymbol{m} + \boldsymbol{\theta}^{(t)}$
13: $\quad$ **until** tired
14: $\quad \overline{\boldsymbol{\theta}} \leftarrow \frac{1}{t} \boldsymbol{m}$
15: $\quad$ **return** $\overline{\boldsymbol{\theta}}$

# Differential Calculus Refresher

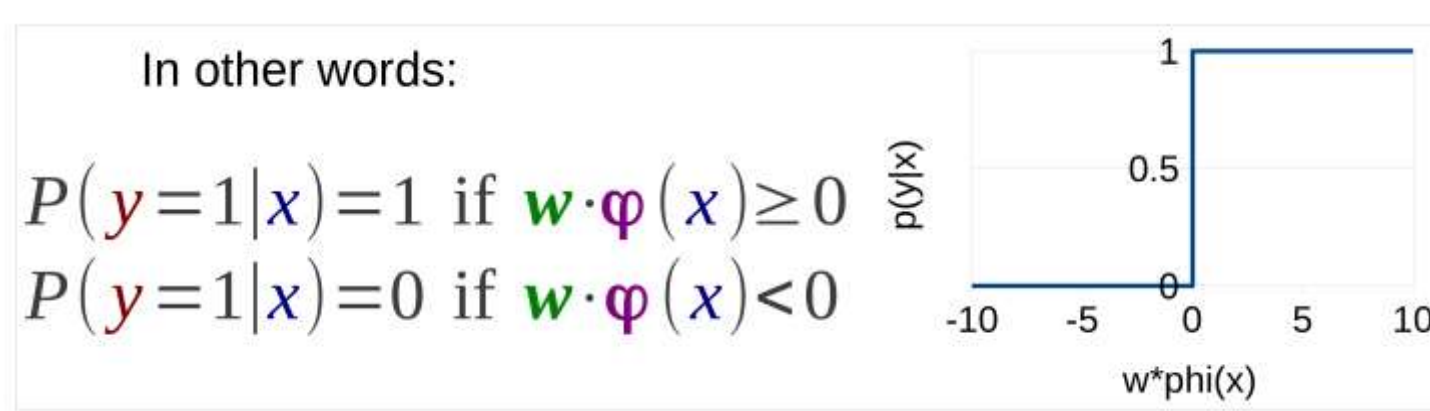- Derivatives
- Chain rule
- Convex functions
- Gradients



Gradient Vectors Shown at Several Points on the
Surface of cos(x) sin(y)

# Logistic Regression
# for **Binary** Classification

# Perceptron & Probabilities

- What if we want a probability p(y|x)?

- The perceptron gives us a prediction y
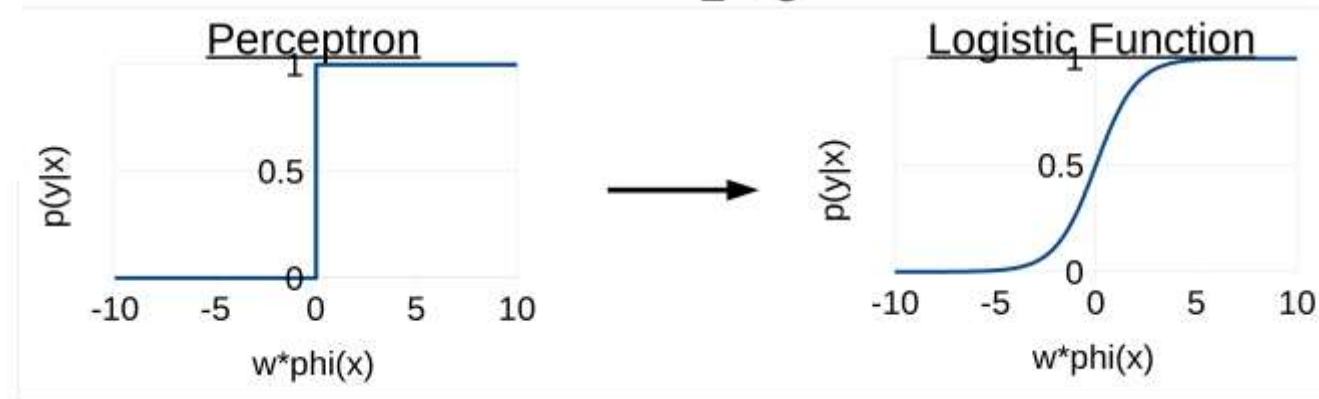  - Let's illustrate this with binary classification

In other words:

$$P(y=1|x)=1 \text{ if } \boldsymbol{w}\cdot\boldsymbol{\varphi}(x)\geq 0$$
$$P(y=1|x)=0 \text{ if } \boldsymbol{w}\cdot\boldsymbol{\varphi}(x)<0$$



Illustrations: Graham Neubig

# The logistic function

$$P(y=1|x)=\frac{e^{w\cdot\varphi(x)}}{1+e^{w\cdot\varphi(x)}}$$

- x: the input
- **φ(x)**: vector of feature functions {$\varphi_1(x)$, $\varphi_2(x)$, …, $\varphi_i(x)$}
- **w**: the weight vector {$w_1$, $w_2$, …, $w_i$}
- y: the prediction, +1 if "yes", -1 if "no"



- "Softer" function than in perceptron
- Can account for uncertainty
- Differentiable

# Logistic regression: how to train?

- Train based on **conditional likelihood**

- Find parameters w that maximize conditional likelihood of all answers $y_i$ given examples $x_i$

$$\hat{w} = \underset{w}{\mathrm{argmax}} \prod_i P(y_i | x_i ; w)$$

# Stochastic gradient ascent (or descent)

- Online training algorithm

```
create map w
for I iterations
    for each labeled pair x, y in the data
        w += α * dP(y|x)/dw
```

- Update weights for every training example
- Move in direction given by gradient
- Size of update step scaled by learning rate

# Gradient of the logistic function

$$\frac{d}{dw}P(y=1|x) = \frac{d}{dw}\frac{e^{w\cdot\varphi(x)}}{1+e^{w\cdot\varphi(x)}}$$

$$= \varphi(x)\frac{e^{w\cdot\varphi(x)}}{(1+e^{w\cdot\varphi(x)})^2}$$

$$\frac{d}{dw}P(y=-1|x) = \frac{d}{dw}(1-\frac{e^{w\cdot\varphi(x)}}{1+e^{w\cdot\varphi(x)}})$$

$$= -\varphi(x)\frac{e^{w\cdot\varphi(x)}}{(1+e^{w\cdot\varphi(x)})^2}$$

# Example: Person/not-person classification problem

Given an introductory sentence in Wikipedia

predict whether the article is about a person

| Given | Predict |
| --- | --- |
| Gonso was a Sanron sect priest (754-827) in the late Nara and early Heian periods. → | Yes! |
| Shichikuzan Chigogataki Fudomyoo is a historical site located at Magura, Maizuru City, Kyoto Prefecture. → | No! |

# Example: initial update

- Set α=1, initialize **w=0**

**x** = A site , located in Maizuru , Kyoto    y = -1

$$\boldsymbol{w} \cdot \boldsymbol{\varphi}(x) = 0 \qquad \frac{d}{dw} P(y=-1|x) = -\frac{e^0}{(1+e^0)^2} \boldsymbol{\varphi}(x)$$

$$= -0.25 \boldsymbol{\varphi}(x)$$

$$\boldsymbol{w} \leftarrow \boldsymbol{w} + -0.25 \boldsymbol{\varphi}(x)$$

| | | | |
|---|---|---|---|
| W unigram "Maizuru" | = -0.25 | W unigram "A" | = -0.25 |
| W unigram "," | = -0.5 | W unigram "site" | = -0.25 |
| W unigram "in" | = -0.25 | W unigram "located" | = -0.25 |
| W unigram "Kyoto" | = -0.25 | | |

# Example: second update

$$\mathbf{x} = \text{Shoken , monk born in Kyoto} \qquad y = 1$$

$$\overset{-0.5 \qquad\qquad\qquad -0.25 \quad -0.25}{\mathbf{w} \cdot \boldsymbol{\varphi}(x) = -1} \qquad \frac{d}{dw} P(y=1|x) \;=\; \frac{e^1}{(1+e^1)^2} \boldsymbol{\varphi}(x)$$

$$= \quad 0.196\, \boldsymbol{\varphi}(x)$$

$$\mathbf{w} \leftarrow \mathbf{w} + 0.196\, \boldsymbol{\varphi}(x)$$

| | | | | | |
|---|---|---|---|---|---|
| $w_{\text{unigram "Maizuru"}}$ | = -0.25 | $w_{\text{unigram "A"}}$ | = -0.25 | $w_{\text{unigram "Shoken"}}$ | = 0.196 |
| $w_{\text{unigram ","}}$ | = -0.304 | $w_{\text{unigram "site"}}$ | = -0.25 | $w_{\text{unigram "monk"}}$ | = 0.196 |
| $w_{\text{unigram "in"}}$ | = -0.054 | $w_{\text{unigram "located"}}$ | = -0.25 | $w_{\text{unigram "born"}}$ | = 0.196 |
| $w_{\text{unigram "Kyoto"}}$ | = -0.054 | | | | |

# How to set the learning rate?

- Various strategies
  - decay over time

$$\alpha = \frac{1}{C + t}$$

Parameter

Number of samples

- Use held-out test set, increase learning rate when likelihood increases

# What you should know about linear models

- Standard supervised learning set-up for text classification
  - Difference between train vs. test data
  - How to evaluate
- 3 examples of linear classifiers
  - Naïve Bayes, Perceptron, Logistic Regression
    - How to make predictions, how to train, strengths and weaknesses
  - Learning as optimization: loss functions and their properties
  - Difference between generative vs. discriminative classifiers
- General machine learning concepts
  - Smoothing, overfitting, underfitting, regularization