

# Linear Models Continued: Perceptron & Logistic Regression

CMSC 723 / LING 723 / INST 725

Marine Carpuat

Slides credit: Graham Neubig, Jacob  
Eisenstein

# Linear Models for Classification

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

Feature  
function  
representation

Weights

# Naïve Bayes recap

- Define  $p(\mathbf{x}, \mathbf{y})$  via a *generative model*
- Prediction:  $\hat{y} = \arg \max_y p(\mathbf{x}_i, y)$
- Learning:

$$\boldsymbol{\theta} = \arg \max_{\boldsymbol{\theta}} p(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta})$$

$$p(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta}) = \prod_i p(\mathbf{x}_i, y_i; \boldsymbol{\theta}) = \prod_i p(\mathbf{x}_i | y_i) p(y_i)$$

$$\phi_{y,j} = \frac{\sum_{i:Y_i=y} x_{ij}}{\sum_{i:Y_i=y} \sum_j x_{ij}}$$

$$\mu_y = \frac{\text{count}(Y = y)}{N}$$

This gives the maximum likelihood estimator (MLE; same as relative frequency estimator)

# The Perceptron

# The perceptron

- A linear model for classification
- An algorithm to learn feature weights given labeled data
  - online algorithm
  - error-driven

# Multiclass perceptron

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

---

**Algorithm 1** Perceptron learning algorithm

---

```
1: procedure PERCEPTRON( $\mathbf{x}_{1:N}, y_{1:N}$ )
2:   repeat
3:     Select an instance  $i$ 
4:      $\hat{y} \leftarrow \arg \max_y \boldsymbol{\theta}_t^\top \mathbf{f}(\mathbf{x}_i, y)$ 
5:     if  $\hat{y} \neq y_i$  then
6:        $\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \mathbf{f}(\mathbf{x}_i, y_i) - \mathbf{f}(\mathbf{x}_i, \hat{y})$ 
7:     else
8:       do nothing
9:   until tired
```

---

# Understanding the perceptron

- What's the impact of the update rule on parameters?
- The perceptron algorithm will converge if the training data is **linearly separable**
  - Proof: see ["A Course In Machine Learning" Ch.4](#)
- Practical issues
  - How to initialize?
  - When to stop?
  - How to order training examples?

# When to stop?

- One technique
  - When the accuracy on held out data starts to decrease
  - Early stopping



Requires splitting data into 3 sets:  
training/development/test



ML fundamentals aside:  
overfitting/underfitting/generalization

# Training error is not sufficient

- We care about **generalization** to new examples
- A classifier can classify training data perfectly, yet classify new examples incorrectly
  - Because training examples are only a sample of data distribution
    - a feature might correlate with class by coincidence
  - Because training examples could be noisy
    - e.g., accident in labeling

# Overfitting

- Consider a model  $\theta$  and its:
  - Error rate over training data  $error_{train}(\theta)$
  - True error rate over all data  $error_{true}(\theta)$
- We say  $h$  overfits the training data if
$$error_{train}(\theta) < error_{true}(\theta)$$

# Evaluating on test data

- Problem: we don't know  $error_{true}(\theta)$ !
- Solution:
  - we set aside a test set
    - some examples that will be used for evaluation
  - we don't look at them during training!
  - after learning a classifier  $\theta$ , we calculate  $error_{test}(\theta)$

# Overfitting

- Another way of putting it
- A classifier  $\theta$  is said to **overfit the training data**, if there is another hypothesis  $\theta'$ , such that
  - $\theta$  has a smaller error than  $\theta'$  on the training data
  - but  $\theta$  has larger error on the test data than  $\theta'$ .

# Underfitting/Overfitting

- Underfitting

- Learning algorithm had the opportunity to learn more from training data, but didn't

- Overfitting

- Learning algorithm paid too much attention to idiosyncracies of the training data; the resulting classifier doesn't generalize

Back to the Perceptron

# Averaged Perceptron improves generalization

---

**Algorithm 2** Averaged perceptron learning algorithm

---

```
1: procedure AVG-PERCEPTRON( $\mathbf{x}_{1:N}, y_{1:N}$ )
2:   repeat
3:     Select an instance  $i$ 
4:      $\hat{y} \leftarrow \arg \max_y \boldsymbol{\theta}_t^\top \mathbf{f}(\mathbf{x}_i, y)$ 
5:     if  $\hat{y} \neq y_i$  then
6:        $\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \mathbf{f}(\mathbf{x}_i, y_i) - \mathbf{f}(\mathbf{x}_i, \hat{y})$ 
7:        $\mathbf{m} \leftarrow \mathbf{m} + \boldsymbol{\theta}_{t+1}$ 
8:     else
9:       do nothing
10:  until tired
11:   $\bar{\boldsymbol{\theta}} \leftarrow \frac{1}{t} \mathbf{m}$ 
```

---



# What objective/loss does the perceptron optimize?

- Zero-one loss function

$$\ell_{\text{perceptron}}(\boldsymbol{\theta}; \mathbf{x}_i, y_i) = \begin{cases} 0, & y_i = \arg \max_y \boldsymbol{\theta}^\top \mathbf{f}(x_i, y) \\ 1, & \text{otherwise} \end{cases}$$

- What are the pros and cons compared to Naïve Bayes loss?

# Logistic Regression

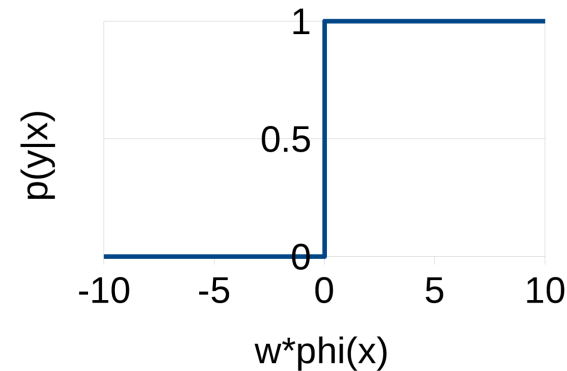
# Perceptron & Probabilities

- 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:

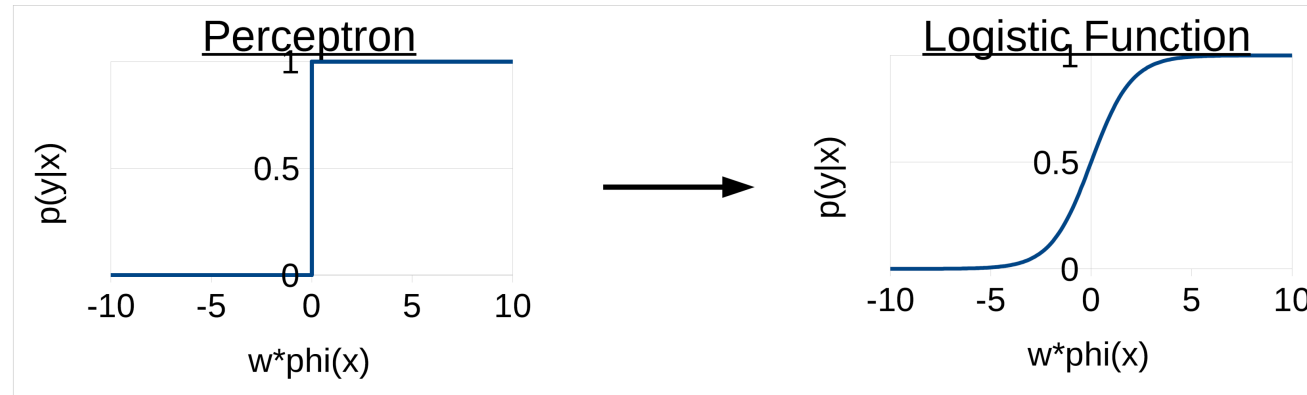
$$P(y=1|x)=1 \text{ if } w \cdot \varphi(x) \geq 0$$

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



# The logistic function

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



- “Softer” function than in perceptron
- Can account for uncertainty
- Differentiable

# Logistic regression: how to train?

- Train based on **conditional likelihood**
- Find parameters  $\mathbf{w}$  that maximize conditional likelihood of all answers  $y_i$  given examples  $x_i$

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

# Stochastic gradient ascent (or descent)

- Online training algorithm for logistic regression
  - and other probabilistic models

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

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

# What you should know

- Standard supervised learning set-up for text classification
  - Difference between train vs. test data
  - How to evaluate
- 3 examples of supervised linear classifiers
  - Naïve Bayes, Perceptron, Logistic Regression
  - Learning as optimization: what is the objective function optimized?
  - Difference between generative vs. discriminative classifiers
  - Smoothing, regularization
  - Overfitting, underfitting

# An online learning algorithm

```
create map  $w$ 
for / iterations
    for each labeled pair  $x, y$  in the data
         $\phi = \text{CREATE\_FEATURES}(x)$ 
         $y' = \text{PREDICT\_ONE}(w, \phi)$ 
        if  $y' \neq y$ 
             $\text{UPDATE\_WEIGHTS}(w, \phi, y)$ 
```



# Perceptron weight update

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

- If  $y = 1$ , increase the weights for features in  $\boldsymbol{\varphi}(\mathbf{x})$
- If  $y = -1$ , decrease the weights for features in  $\boldsymbol{\varphi}(\mathbf{x})$