

## Perceptron learning algorithm

<https://powcoder.com>

---

The multiclass perceptron learning algorithm

---

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

---

---

## The averaged perceptron algorithm

---

```
1: procedure AVERAGE-PERCEPTRON( $\mathbf{x}^{1:N}, \mathbf{y}^{1:N}$ )
2:    $t \leftarrow 0$ 
3:    $\boldsymbol{\theta}^{(0)} \leftarrow \mathbf{0}$ 
4:    $\mathbf{m} \leftarrow \mathbf{0}$ 
5:   repeat
6:      $t \leftarrow t + 1$ 
7:     Select an instance  $i$  ▷ Online training
8:      $\hat{y} \leftarrow \arg\max_y \boldsymbol{\theta}^{(t-1)} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y)$ 
9:     if  $\hat{y} \neq y^{(i)}$  then
10:       $\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)} + \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \mathbf{f}(\mathbf{x}^{(i)}, \hat{y})$ 
11:    else
12:       $\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)}$ 
13:    end if
14:     $\mathbf{m} \leftarrow \mathbf{m} + \boldsymbol{\theta}^{(t)}$ 
15:  until tired
16:   $\bar{\boldsymbol{\theta}} = \frac{1}{t} \mathbf{m}$ 
17:  return  $\bar{\boldsymbol{\theta}}$ 
18: end procedure
```

Making predictions

<https://powcoder.com>

## Assignment Project Exam Help

Assuming that you have trained weights  $\theta$ , finding the best scoring label is straightforward given a data instance  $x^{(i)}$ :

<https://powcoder.com>

$$\hat{y} = \operatorname{argmax}_y \theta \cdot f(x^{(i)}, y)$$

Add WeChat powcoder

Note that the score cannot be interpreted probabilistically.

What is the Loss function of Perceptron?

<https://powcoder.com>

## Assignment Project Exam Help

- ▶ Hinge Loss:

$$\ell_{\text{perceptron}}(\theta, \mathbf{x}^{(i)}, y^{(i)}) = \max_{y \in \mathcal{Y}} \theta \cdot \mathbf{f}(\mathbf{x}^{(i)}, y) - \theta \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)})$$

- ▶ When  $\hat{y} = y^{(i)}$ , the loss is zero, otherwise, it increases linearly with the gap between the score for the predicted label  $\hat{y}$  and the score for the true label  $y^{(i)}$ .
- ▶ Plot the loss against the input gives a Hinge shape, giving its name **Hinge loss**

## What is the Loss function of Perceptron?

<https://powcoder.com>

- Hinge Loss:

$$\ell_{\text{perceptron}}(\theta; \mathbf{x}^{(i)}, y^{(i)}) = \max_{y \in \mathcal{Y}} \theta \cdot \mathbf{f}(\mathbf{x}^{(i)}, y) - \theta \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)})$$

- The derivative of the loss function:

$$\frac{\partial}{\partial \theta} \ell_{\text{perceptron}}(\theta; \mathbf{x}^{(i)}, y^{(i)}) = \mathbf{f}(\mathbf{x}^{(i)}, \hat{y}) - \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)})$$

- Updating: At each instance, the perceptron algorithm takes a step of magnitude one in the opposite direction of this gradient

$$\begin{aligned} \theta &= \theta - \nabla_{\theta} \ell_{\text{perceptron}} = \theta - \frac{\partial}{\partial \theta} \ell_{\text{perceptron}}(\theta; \mathbf{x}^{(i)}, y^{(i)}) \\ &= \theta - \mathbf{f}(\mathbf{x}^{(i)}, \hat{y}) + \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) \end{aligned}$$

When do you stop the training?

<https://powcoder.com>

How many iterations do you have to train, and when do you stop?

- ▶ Ideally you want to stop where the model have the best performance on previously unseen data.
- ▶ One way is to check the difference between the average weight vectors after each pass of the data. If the norm of the difference falls below a certain predefined threshold, then stop training.
- ▶ **Early stopping.** Hold out a proportion of the training data, and when the accuracy on this held out data set starts to decrease, the model has begun to **overfit** the training set. It's time to stop training.



## Multiclass Perceptron: Parameter estimation

<https://powcoder.com>

Assignment Project Exam Help

	not	funny	painful	ok	overall	story	good	jokes	<i>bias</i>
POS	0	0	0	0	0	0	0	0	1
NEG	0	0	0	0	0	0	0	0	0
NEU	0	0	0	0	0	0	0	0	0

Training instance:  $y = \text{NEG}$ ,  $x = \text{"not funny at all"}$

Add WeChat powcoder

$$\hat{y} = \underset{y}{\operatorname{argmax}} \text{score}(y, \mathbf{x}) = \text{POS}$$

$y \neq y'$ , so update



## Multi-class Perceptron: parameter estimation

<https://powcoder.com>

### Assignment Project Exam Help

Updated  $\theta$

	not	funny	painful	etc	overall	story	good	books	
POS	-1	-1	0	0	0	0	0	0	0
NEG	1	1	0	0	0	0	0	0	1
NEU	0	0	0	0	0	0	0	0	0

Updating:

Add WeChat powcoder

- ▶ Add the feature vector of the instance  $x^{(i)}$  to weights for the *NEG* class and subtract from the weights for the *POS* class.
- ▶ The weights for *NEU* are not updated. Why not?

## Multi-class Perceptron: parameter estimation

<https://powcoder.com>

	not	funny	painful	ok	overall	story	good	ickes	bias
POS	-1	-1	0	0	0	0	0	0	0
NEG	1	1	0	0	0	0	0	0	1
NEU	0	0	0	0	0	0	0	0	0

New Training instance:  $y = \text{NEG}$   $x = \text{"painful, not funny"}$

$$\text{score}(y, x) = \theta \cdot f$$

$$\text{score}(y = \text{POS}, x) = 1 \times -1 + 1 \times -1 + 1 \times 0 + 0 \times 0 + 0 \times 0 + 0 \times 0 + 0 \times 0 + 0 \times 0 + 1 \times 0 = -2$$

$$\text{score}(y = \text{NEG}, x) = 1 \times 1 + 1 \times 1 + 1 \times 0 + 0 \times 0 + 0 \times 0 + 0 \times 0 + 0 \times 0 + 0 \times 0 + 1 \times 1 = 3$$

$$\text{score}(y = \text{NEU}, x) = 1 \times 0 + 1 \times 0 + 1 \times 0 + 0 \times 0 + 0 \times 0 + 0 \times 0 + 0 \times 0 + 0 \times 0 + 1 \times 0 = 0$$

$\hat{y} = \text{NEG}$ , correct, so no update

## Perceptron versus Naïve Bayes

- ▶ Both  $\ell_{NB}$  and  $\ell_{PERCEPTRON}$  are convex, making them relatively easy to optimize.  $\ell_{NB}$  can be optimized in closed form, while  $\ell_{PERCEPTRON}$  requires iterating over the training set multiple times.
- ▶  $\ell_{NB}$  can suffer infinite loss on a single example, since the logarithm of zero probability is negative infinity. Naïve Bayes will therefore overemphasize some examples and underemphasize others.
- ▶ The Naïve Bayes classifier assumes that the observed features are conditionally independent, given the label, and the performance of the classifier depends on the extent to which this assumption holds. The perceptron requires no such assumption.
- ▶  $\ell_{PERCEPTRON}$  treats all correct answers equally. Even if  $\theta$  only gives the correct answer by a tiny margin, the loss is still zero.

## Perceptron vs Logistic regression

<https://powcoder.com>

- ▶ Inference: Both Perceptron and Logistic Regression are discriminative models, but the score for Logistic Regression can be interpreted probabilistically, while the score for Perceptron can not.
- ▶ Parameter Estimation. The parameters for both Logistic Regression and Perceptron are estimated iteratively by gradient descent. The difference in parameter estimation between Logistic Regression and Perceptron follows from their loss functions.
  - ▶ The features for which the weights need to be updated for Perceptron are features for the correct label and the label that receives the highest score if they are different. For Logistic Regression, the features “fire” for all labels are updated.
  - ▶ In perceptron, incorrect features penalized by 1, while for Logistic Regression incorrect features penalized proportionally by how much off is its prediction is.

## Online large margin classification

<https://powcoder.com>

- ▶ For the perceptron, if the classifier gets the correct answer on the training example by a small margin, it may give a different answer on a nearby test instance.
- ▶ The address this we need the answer to be correct by a large margin. **Margin** can be formalized as:

$$\gamma(\theta; \mathbf{x}^{(i)}, y^{(i)}) = \theta \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \max_{y \neq y^{(i)}} \theta \cdot \mathbf{f}(\mathbf{x}^{(i)}, y)$$

- ▶ The margin represents the difference between the score of the correct label and the score for the highest-scoring incorrect label.
- ▶ The intuition is that it's not enough to label the training data correctly — the correct label should be separated from other labels by a comfortable margin.

## Margin Loss

<https://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help

$$\ell_{\text{MARGIN}}(\theta; \mathbf{x}^{(i)}, y^{(i)}) = \begin{cases} 1 - \gamma(\theta; \mathbf{x}^{(i)}, y^{(i)}), & \text{otherwise} \\ 0, & \gamma(\theta; \mathbf{x}^{(i)}, y^{(i)}) \geq 1 \end{cases}$$

<https://powcoder.com>

$$= (1 - \gamma(\theta; \mathbf{x}^{(i)}, y^{(i)}))_+$$

Add WeChat powcoder

where  $(x)_+ = \max(0, x)$ . The margin loss is zero when the margin between the score for the true label and the best alternative  $\hat{y}$  is at least 1.

## Minimizing the margin loss

The margin can be achieved by adding a cost to the score

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} \theta \cdot \mathbf{f}(\mathbf{x}^{(i)}, y) + c(y^{(i)}, y)$$

Assignment Project Exam Help

where the cost function that can be defined as

$$c(y^{(i)}, y) = \begin{cases} 1, & \text{if } y \neq y^{(i)} \\ 0, & \text{otherwise} \end{cases}$$

Note:

- ▶ The cost function is only used to compute  $\hat{y}$  during training as the true label is unknown for unseen data. During test time the best label is just the label with the best score.
- ▶ For purposes of training, don't simply use the label that has the highest score. Instead, choose the one that has the highest score and cost.
- ▶ When the label that has the highest score, it may still be chosen when it has a score that is higher than an (incorrect) alternative label by at least 1, in which case the loss will be 0

The update rule for SVM

<https://powcoder.com>

Finding the gradient of the loss function involves constrained optimization with Lagrangian multipliers. The final update rule is:

[Add WeChat powcoder](https://powcoder.com)

$$\theta^{(t)} \leftarrow (1 - \lambda)\theta^{(t-1)} + \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \mathbf{f}(\mathbf{x}^{(i)}, \hat{y})$$

<https://powcoder.com>

- ▶ The previous weights  $\theta^{(t-1)}$  are scaled by  $1 - \lambda$ , where  $\lambda \in \{0, 1\}$ . [Add WeChat powcoder](https://powcoder.com)
- ▶ It pulls the weights back towards zero, and in this sense it serves as a form of regularization that prevents overfitting.
- ▶ How is this different from  $L2$  regularization?



## Linear algorithms: a summary

- ▶ **Naïve Bayes:** *Pros:* easy to implement, estimation is fast, required on a single pass on the training data; assigns probabilities to predicted labels; controls overfitting with parameter smoothing; *Cons:* often has poor accuracy, especially with correlated features
- ▶ **Perceptron:** *Pros:* easy to implement, online, error-driven learning typically produces good accuracy, especially after averaging. *Cons:* not probabilistic, hard to know when to stop learning; lack of margin can lead to overfitting.
- ▶ **Support vector machine:** *Pros:* optimizes an error-based metric, usually resulting in good accuracy; overfitting is controlled by a regularization parameter. *Cons:* not probabilistic.
- ▶ **Logistic regression:** *Pros:* error-driven and probabilistic; overfitting is controlled by a regularization parameter. *Cons:* batch learning requires black-box optimization; Logistic loss can “overtrain” on correctly labeled samples