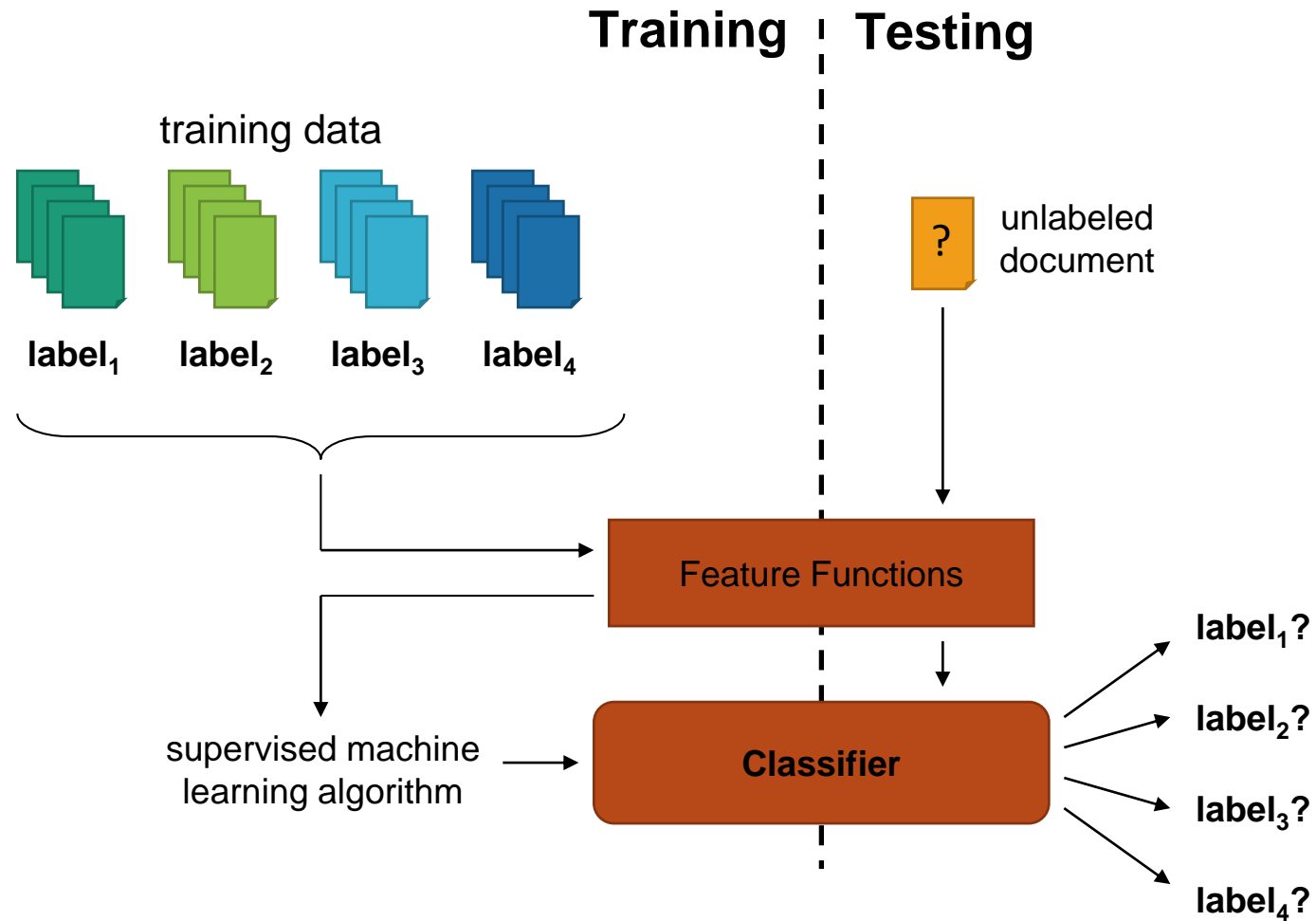# Supervised Classification with the Perceptron

**CMSC 470**

Marine Carpuat

# Last time

- Word senses distinguish different meanings of same word
- Sense inventories
- Annotation issues and annotator agreement (Kappa)
- Definition of Word Sense Disambiguation Task
- An unsupervised approach: Lesk algorithm
- Supervised classification:
  - Train vs. test data
  - The most frequent class baseline
- Evaluation metrics: accuracy, precision, recall

# WSD as **Supervised** Classification

# Evaluation Metrics for Classification

# How are annotated examples used in supervised learning?

- Supervised learning = requires examples annotated with correct prediction

- Used in 2 ways:
  - To find good values for the model (hyper)parameters (training data)
  - To evaluate how good the resulting classifier is (test data)

- How do we know how good a classifier is?
  - Compare classifier predictions with human annotation
  - On held out test examples
  - Evaluation metrics: accuracy, precision, recall

# Quantifying Errors in a Classification Task: The 2-by-2 contingency table (per class)

|  | correct | not correct |
|---|---|---|
| selected | tp | fp |
| not selected | fn | tn |

# Quantifying Errors in a Classification Task: Precision and Recall

|  | correct | not correct |
|---|---|---|
| selected | tp | fp |
| not selected | fn | tn |

**Precision**: % of selected items that are correct
**Recall**: % of correct items that are selected

Q: When are Precision/Recall more informative than accuracy?

# A combined measure: F

- A combined measure that assesses the P/R tradeoff is F measure (weighted harmonic mean):

$$F = \cfrac{1}{a\cfrac{1}{P} + (1-a)\cfrac{1}{R}} = \frac{(b^2+1)PR}{b^2P+R}$$
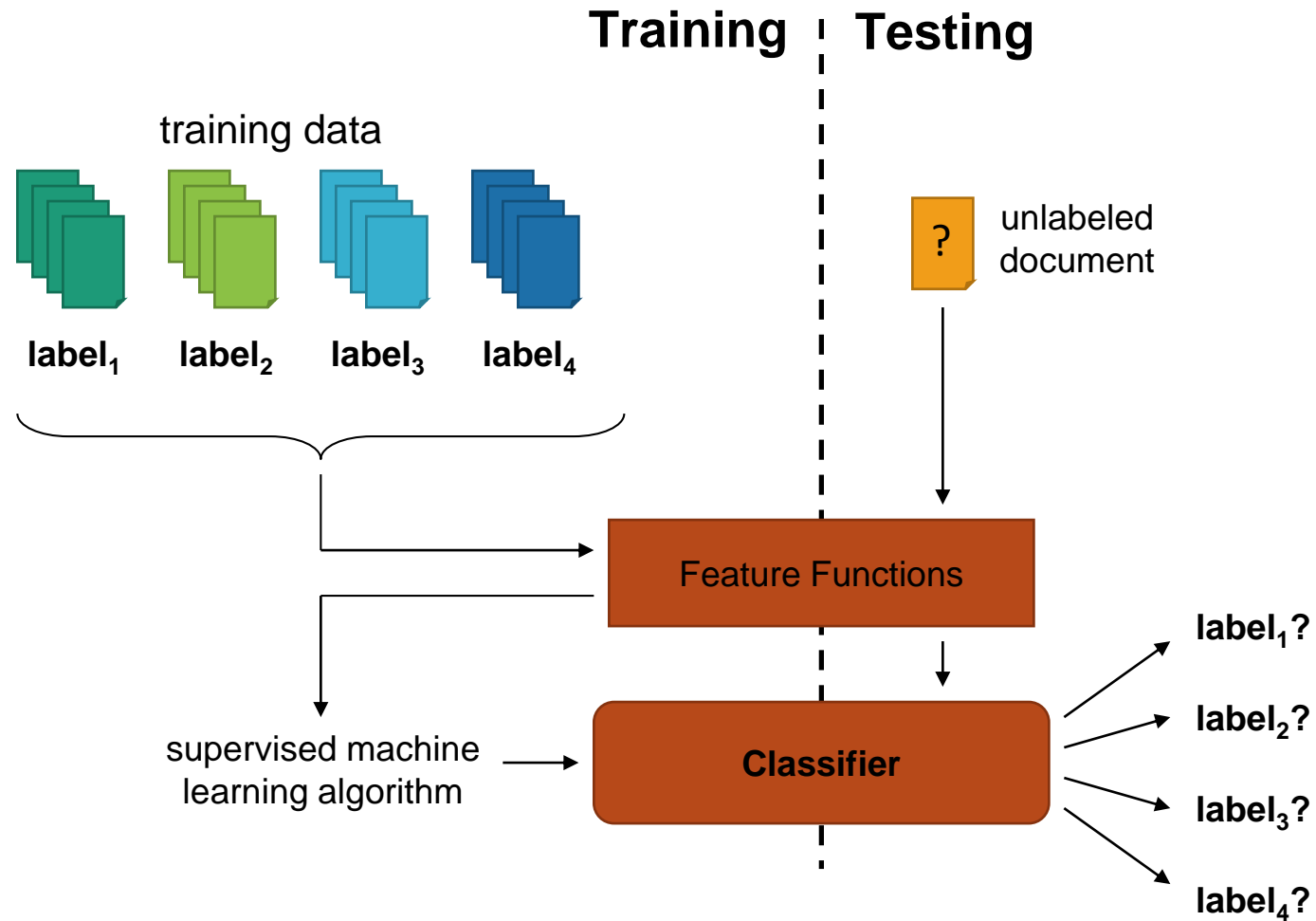
With $\beta^2 = \dfrac{1}{\alpha} - 1$

- People usually use balanced F1 measure
  - i.e., with $\beta = 1$ (that is, $\alpha$ = ½):

    *F = 2PR/(P+R)*

# The Perceptron

A simple Supervised Classifier

# WSD as **Supervised** Classification



**Training** ┊ **Testing**

training data

label₁  label₂  label₃  label₄

? unlabeled document

Feature Functions

supervised machine learning algorithm

**Classifier**

label₁?

label₂?

label₃?

label₄?

# Formalizing classification

**Task definition**

- *Given inputs*:
  - an example $x$

    often x is a D-dimensional vector of binary or real values
  - a fixed set of classes $Y$

    $Y = \{y_1, y_2, ..., y_J\}$

    e.g. word senses from WordNet

- *Output*: a predicted class $y \in Y$

**Classifier definition**

A function $f: x \rightarrow f(x) = y$

Many different types of functions/classifiers can be defined

- We'll talk about perceptron, logistic regression, neural networks.

# Example:
# Word Sense Disambiguation for "bass"

- Y = {-1,+1} since there are 2 senses in our inventory

- Many different definitions of x are possible
  - E.g., vector of word frequencies for words that co-occur in a window of +/- k words around "bass"
  - Instead of frequency, we could use binary values, or tf.idf, or PPMI, etc.
  - Instead of window, we could use the entire sentence
  - Instead of/in addition to words, we could use POS tags
  - …

| WordNet Sense | Spanish Translation | Roget Category | Target Word in Context |
|---|---|---|---|
| bass⁴ | lubina | FISH/INSECT | …fish as Pacific salmon and striped **bass** and… |
| bass⁴ | lubina | FISH/INSECT | …produce filets of smoked **bass** or sturgeon… |
| bass⁷ | bajo | MUSIC | …exciting jazz **bass** player since Ray Brown… |
| bass⁷ | bajo | MUSIC | …play **bass** because he doesn't have to solo… |

# Perception Test Algorithm for Binary Classification: Predict class -1 or +1 for example x

$$f(x) = sign(w.x + b)$$

**Algorithm 6** PERCEPTRONTEST($w_0, w_1, \ldots, w_D, b, \hat{x}$)

1:  $a \leftarrow \sum_{d=1}^{D} w_d \, \hat{x}_d + b$          // compute activation for the test example

2:  **return** SIGN($a$)

# Perceptron Training Algorithm:
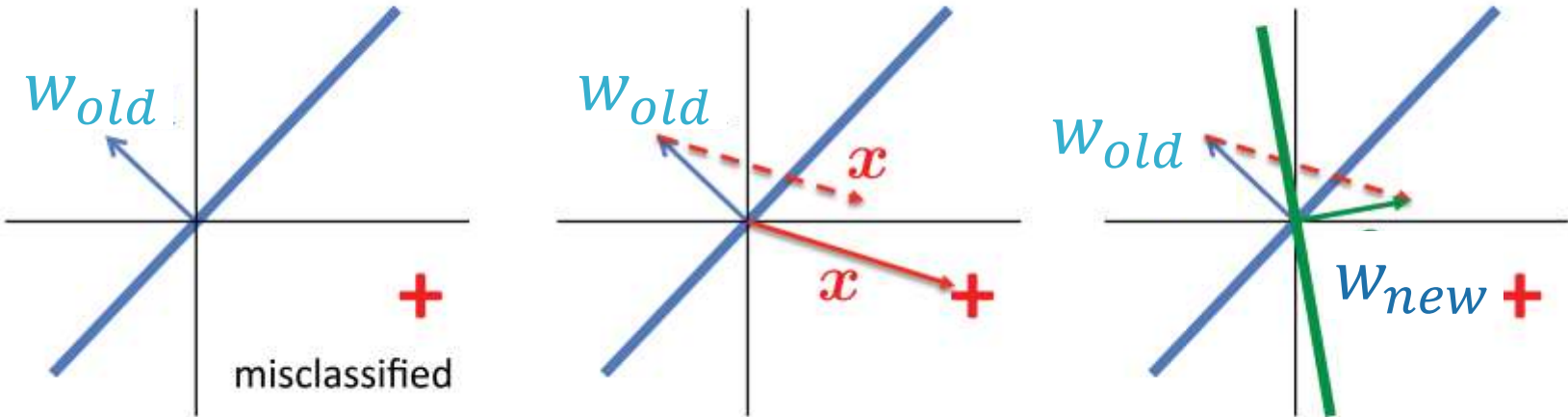# Find good values for (w,b) given training data D

**Algorithm 5** PERCEPTRONTRAIN($\mathbf{D}$, *MaxIter*)

1: $w_d \leftarrow 0$, for all $d = 1 \dots D$     // initialize weights
2: $b \leftarrow 0$     // initialize bias
3: **for** *iter* $= 1 \dots MaxIter$ **do**
4:     **for all** $(x,y) \in \mathbf{D}$ **do**
5:        $a \leftarrow \sum_{d=1}^{D} w_d \, x_d + b$     // compute activation for this example
6:        **if** $ya \leq 0$ **then**
7:           $w_d \leftarrow w_d + yx_d$, for all $d = 1 \dots D$     // update weights
8:           $b \leftarrow b + y$     // update bias
9:        **end if**
10:     **end for**
11: **end for**
12: **return** $w_0, w_1, \dots, w_D, b$

# The Perceptron update rule: geometric interpretation

# Machine Learning Vocabulary

x is often called the feature vector
- its elements are defined (by us, the model designers) to capture properties or features of the input that are expected to correlate with predictions

w and b are the parameters of the classifier
- they are needed to fully define the classification function $f(x) = y$
- their values are found by the training algorithm using training data D

MaxIter is a hyperparameter
- controls when training stops
- MaxIter impacts the nature of function f indirectly

All of the above affect the performance of the final classifier!

# Standard Perceptron: predict based on final parameters

**Algorithm 5** PERCEPTRONTRAIN(**D**, *MaxIter*)

1: $w_d \leftarrow 0$, for all $d = 1 \ldots D$                    // initialize weights
2: $b \leftarrow 0$                    // initialize bias
3: **for** *iter* = 1 ... *MaxIter* **do**
4:     **for all** $(x,y) \in$ **D do**
5:         $a \leftarrow \sum_{d=1}^{D} w_d \, x_d + b$            // compute activation for this example
6:         **if** $ya \leq 0$ **then**
7:             $w_d \leftarrow w_d + yx_d$, for all $d = 1 \ldots D$            // update weights
8:             $b \leftarrow b + y$                    // update bias
9:         **end if**
10:     **end for**
11: **end for**
12: **return** $w_0, w_1, \ldots, w_D, b$

# Predict based on final + intermediate parameters

- The voted perceptron

$$\hat{y} = \text{sign}\left(\sum_{k=1}^{K} c^{(k)}\text{sign}\left(\boldsymbol{w}^{(k)} \cdot \hat{\boldsymbol{x}} + b^{(k)}\right)\right)$$

- The averaged perceptron

$$\hat{y} = \text{sign}\left(\sum_{k=1}^{K} c^{(k)}\left(\boldsymbol{w}^{(k)} \cdot \hat{\boldsymbol{x}} + b^{(k)}\right)\right)$$

- Require keeping track of "survival time" of weight vectors

$$c^{(1)}, \ldots, c^{(K)}$$

# How would you modify this algorithm for voted perceptron?

**Algorithm 5** PERCEPTRONTRAIN(**D**, *MaxIter*)

1: $w_d \leftarrow 0$, for all $d = 1 \ldots D$      // initialize weights
2: $b \leftarrow 0$      // initialize bias
3: **for** $iter = 1 \ldots MaxIter$ **do**
4:      **for all** $(x,y) \in$ **D do**
5:          $a \leftarrow \sum_{d=1}^{D} w_d \, x_d + b$      // compute activation for this example
6:          **if** $ya \leq 0$ **then**
7:              $w_d \leftarrow w_d + yx_d$, for all $d = 1 \ldots D$      // update weights
8:              $b \leftarrow b + y$      // update bias
9:          **end if**
10:      **end for**
11: **end for**
12: **return** $w_0, w_1, \ldots, w_D, b$

# How would you modify this algorithm for averaged perceptron?

**Algorithm 5** PERCEPTRON TRAIN($D$, *MaxIter*)

1: $w_d \leftarrow 0$, for all $d = 1 \ldots D$  // initialize weights
2: $b \leftarrow 0$  // initialize bias
3: **for** $iter = 1 \ldots MaxIter$ **do**
4:   **for all** $(x,y) \in D$ **do**
5:     $a \leftarrow \sum_{d=1}^{D} w_d \, x_d + b$  // compute activation for this example
6:     **if** $ya \leq 0$ **then**
7:       $w_d \leftarrow w_d + yx_d$, for all $d = 1 \ldots D$  // update weights
8:       $b \leftarrow b + y$  // update bias
9:     **end if**
10:   **end for**
11: **end for**
12: **return** $w_0, w_1, \ldots, w_D, b$

# Averaged perceptron decision rule

$$\hat{y} = \text{sign}\left(\sum_{k=1}^{K} c^{(k)}\left(\boldsymbol{w}^{(k)} \cdot \hat{\boldsymbol{x}} + b^{(k)}\right)\right)$$

can be rewritten as

$$\hat{y} = \text{sign}\left(\left(\sum_{k=1}^{K} c^{(k)}\boldsymbol{w}^{(k)}\right) \cdot \hat{\boldsymbol{x}} + \sum_{k=1}^{K} c^{(k)}b^{(k)}\right)$$

# An Efficient Algorithm for Averaged Perceptron Training

**Algorithm 7** AVERAGEDPERCEPTRONTRAIN($\mathbf{D}$, *MaxIter*)

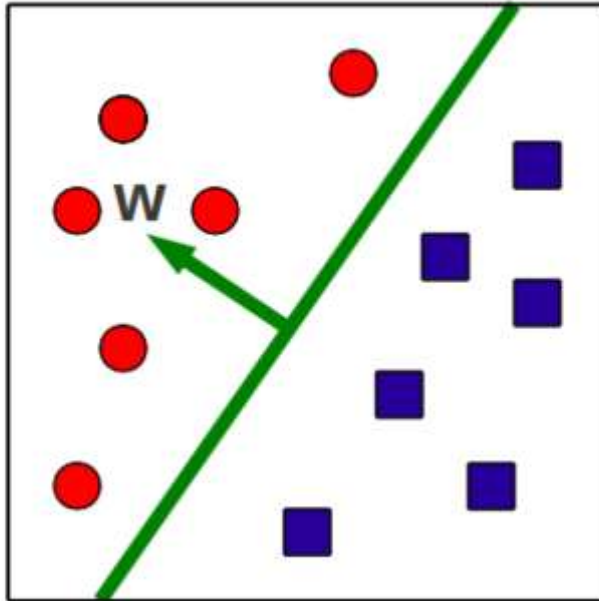| | | |
|---|---|---|
| 1: | $w \leftarrow \langle 0,0,\ldots 0 \rangle \quad , \quad b \leftarrow 0$ | // initialize weights and bias |
| 2: | $u \leftarrow \langle 0,0,\ldots 0 \rangle \quad , \quad \beta \leftarrow 0$ | // initialize cached weights and bias |
| 3: | $c \leftarrow 1$ | // initialize example counter to one |
| 4: | **for** *iter* $= 1 \ldots MaxIter$ **do** | |
| 5: |    **for all** $(x,y) \in \mathbf{D}$ **do** | |
| 6: |       **if** $y(w \cdot x + b) \leq 0$ **then** | |
| 7: |          $w \leftarrow w + y\,x$ | // update weights |
| 8: |          $b \leftarrow b + y$ | // update bias |
| 9: |          $u \leftarrow u + y\,c\,x$ | // update cached weights |
| 10: |          $\beta \leftarrow \beta + y\,c$ | // update cached bias |
| 11: |       **end if** | |
| 12: |       $c \leftarrow c + 1$ | // increment counter regardless of update |
| 13: |    **end for** | |
| 14: | **end for** | |
| 15: | **return** $w - \frac{1}{c} u, \; b - \frac{1}{c} \beta$ | // return averaged weights and bias |

# Perceptron for binary classification



- Classifier = a hyperplane that separates positive from negative examples

$$\hat{y} = sign(w.x + b)$$

- Perceptron training
  - Finds such a hyperplane
  - If training examples are separable

# Convergence of Perceptron

**Theorem (Block & Novikoff, 1962)**

If the training data $D = \{(x_1, y_1), \ldots, (x_N, y_N)\}$ is **linearly separable** with margin $\gamma$ by a unit norm hyperplane $w_*$ ($\|w_*\| = 1$) with $b = 0$,

Then **perceptron training converges after** $\dfrac{R^2}{\gamma^2}$ **errors** during training

(assuming ($\|x\| < R$) for all $x$).

# More Machine Learning vocabulary: 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 are other parameters $\theta'$, 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

- Practical strategies to improve generalization for the perceptron

  - Voting/Averaging

  - Randomize order of training data

  - Use a development test set to find good hyperparameter values
    - E.g., early stopping is a good strategy to avoid overfitting

# The Perceptron
## What you should know

- What is the underlying function used to make predictions
- Perceptron test algorithm
- Perceptron training algorithm
- How to improve perceptron training with the averaged perceptron
- Fundamental Machine Learning Concepts:
  - train vs. test data; parameter; hyperparameter; generalization; overfitting; underfitting.
- How to define features