

# Feature-based Classification with the Perceptron

Nathan Schneider  
(some slides borrowed from Chris Dyer)  
ANLP | 11 October 2017

# Feature-Based Classification

# Word Sense Disambiguation (WSD)

- Given a word **in context**, predict which sense is being used.
  - Evaluated on corpora such as **SemCor**, which is fully annotated for WordNet synsets.
- For example: consider joint POS & WSD classification for ‘interest’, with 3 senses:
  - **N:financial** (*I repaid the loan with **interest***)
  - **N:nonfinancial** (*I read the news with **interest***)
  - **V:nonfinancial** (*Can I **interest** you in a dessert?*)

# Beyond BoW

- Neighboring words are relevant to this decision.
- More generally, we can define **features** of the input that may help identify the correct class.
  - Individual words
  - Bigrams (pairs of consecutive words: *Wall Street*)
  - Capitalization (*interest* vs. *Interest* vs. *INTEREST*)
  - Metadata: document genre, author, ...
- These can be used in naïve Bayes—“bag of features”
  - With overlapping features, independence assumption is *even more naïve*:  $p(y \mid \mathbf{x}) \propto p(y) \cdots p(\text{Wall} \mid y) p(\text{Street} \mid y) p(\text{Wall Street} \mid y)$
  - But other kinds of feature-based classifiers don't make this naïve assumption.

# Feature Extraction

$x$  = Wall Street vets raise concerns  
about **interest** rates , politics

## spelling feature

	$\phi(x)$
bias	1
capitalized?	0
#wordsBefore	6
#wordsAfter	3
relativeOffset	0.66
leftWord=about	1
leftWord=best	0
rightWord=rates	1
rightWord=in	0
Wall	1
Street	1
vets	1
best	0
in	0
Wall Street	1
Street vets	1
vets raise	1

...

- Turns the input into a table of features with real values (often binary: 0 or 1).
- In practice: define feature templates like “leftWord=•” from which specific features are instantiated

# Feature Extraction

$x$  = Wall Street vets raise concerns about **interest** rates , politics

	$\phi(x)$
bias	1
capitalized?	0
#wordsBefore	6
#wordsAfter	3
relativeOffset	0.66
leftWord=about	1
leftWord=best	0
rightWord=rates	1
rightWord=in	0
Wall	1
Street	1
vets	1
best	0
in	0
Wall Street	1
Street vets	1
vets raise	1

## token positional features

- Turns the input into a table of features with real values (often binary: 0 or 1).
- In practice: define feature templates like “leftWord=•” from which specific features are instantiated

...

# Feature Extraction

$x$  = Wall Street vets raise concerns about **interest** rates , politics

	$\phi(x)$
bias	1
capitalized?	0
#wordsBefore	6
#wordsAfter	3
relativeOffset	0.66
leftWord=about	1
leftWord=best	0
rightWord=rates	1
rightWord=in	0
Wall	1
Street	1
vets	1
best	0
in	0
Wall Street	1
Street vets	1
vets raise	1

...

## immediately neighboring words

- Turns the input into a table of features with real values (often binary: 0 or 1).
- In practice: define feature templates like “leftWord=•” from which specific features are instantiated

# Feature Extraction

$x$  = Wall Street vets raise concerns about **interest** rates , politics

	$\phi(x)$
bias	1
capitalized?	0
#wordsBefore	6
#wordsAfter	3
relativeOffset	0.66
leftWord=about	1
leftWord=best	0
rightWord=rates	1
rightWord=in	0
Wall	1
Street	1
vets	1
best	0
in	0
Wall Street	1
Street vets	1
vets raise	1

**unigrams**

- Turns the input into a table of features with real values (often binary: 0 or 1).
- In practice: define feature templates like “leftWord=•” from which specific features are instantiated

...



# Feature Extraction

$x$  = Wall Street vets raise concerns about **interest** rates , politics

	$\phi(x)$
bias	1
capitalized?	0
#wordsBefore	6
#wordsAfter	3
relativeOffset	0.66
leftWord=about	1
leftWord=best	0
rightWord=rates	1
rightWord=in	0
Wall	1
Street	1
vets	1
best	0
in	0
Wall Street	1
Street vets	1
vets raise	1

**bigrams**

...

- Turns the input into a table of features with real values (often binary: 0 or 1).
- In practice: define feature templates like “leftWord=•” from which specific features are instantiated

# Feature Extraction

$x$  = Wall Street vets raise concerns about **interest** rates , politics

	$\phi(x)$
bias	1
capitalized?	0
#wordsBefore	6
#wordsAfter	3
relativeOffset	0.66
leftWord=about	1
leftWord=best	0
rightWord=rates	1
rightWord=in	0
Wall	1
Street	1
vets	1
best	0
in	0
Wall Street	1
Street vets	1
vets raise	1

...

**bias feature** ( $\approx$ class prior): value of 1 for every  $x$  so the learned weight will reflect prevalence of the class

- Turns the input into a table of features with real values (often binary: 0 or 1).
- In practice: define feature templates like “leftWord=•” from which specific features are instantiated

# Feature Extraction

$x$  = Wall Street vets raise concerns about **interest** rates , politics

$x'$  = Pet 's best **interest** in mind , but vets must follow law

	$\phi(x)$	$\phi(x')$
bias	1	1
capitalized?	0	0
#wordsBefore	6	3
#wordsAfter	3	8
relativeOffset	0.66	0.27
leftWord=about	1	0
leftWord=best	0	1
rightWord=rates	1	0
rightWord=in	0	1
Wall	1	0
Street	1	0
vets	1	1
best	0	1
in	0	1
Wall Street	1	0
Street vets	1	0
vets raise	1	0

...

- Turns the input into a table of features with real values (often binary: 0 or 1).
- In practice: define feature templates like “leftWord=•” from which specific features are instantiated

# Choosing Features

- Supervision means that we don't have to pre-specify the precise relationship between each feature and the classification outcomes.
- But domain expertise helps in choosing which kinds of features to include in the model. (words, subword units, metadata, ...)
  - And sometimes, highly task-specific features are helpful.
- The decision about what features to include in a model is called **feature engineering**.
  - (There are some algorithmic techniques, such as *feature selection*, that can assist in this process.)
  - More features = more flexibility, but also more expensive to train, more opportunity for overfitting.

# Linear Model

- For each input  $\mathbf{x}$  (e.g., a document or word token), let  $\phi(\mathbf{x})$  be a function that extracts a vector of its features.
  - Features may be binary (e.g., capitalized?) or real-valued (e.g., #word=debt).
- Each feature receives a real-valued **weight** parameter  $w$ . Each candidate label  $y'$  is scored for the token by summing the weights for the active features:

$$\begin{aligned} & \mathbf{w}_{y'}^\top \phi(\mathbf{x}) \\ &= \sum_j w_{y',j} \cdot \phi_j(\mathbf{x}) \end{aligned}$$

- For binary classification, equivalent to:  $\text{sign}(\mathbf{w}^\top \phi(\mathbf{x}))$  — **+1** or **-1**

	$\phi(\mathbf{x})$	w	$\phi(\mathbf{x}')$
bias	1	-3.00	1
capitalized?	0	.22	0
#wordsBefore	6	-.01	3
#wordsAfter	3	.01	8
relativeOffset	0.6	1.00	0.2
leftWord=about	1	.00	0
leftWord=best	0	-2.00	1
rightWord=rates	1	5.00	0
rightWord=in	0	-1.00	1
Wall	1	1.00	0
Street	1	-1.00	0
vets	1	-.05	1
best	0	-1.00	1
in	0	-.01	1
Wall Street	1	4.00	0
Street vets	1	.00	0
vets raise	1	.00	0

...

$\mathbf{x}$  = Wall Street vets raise concerns about **interest** rates , politics

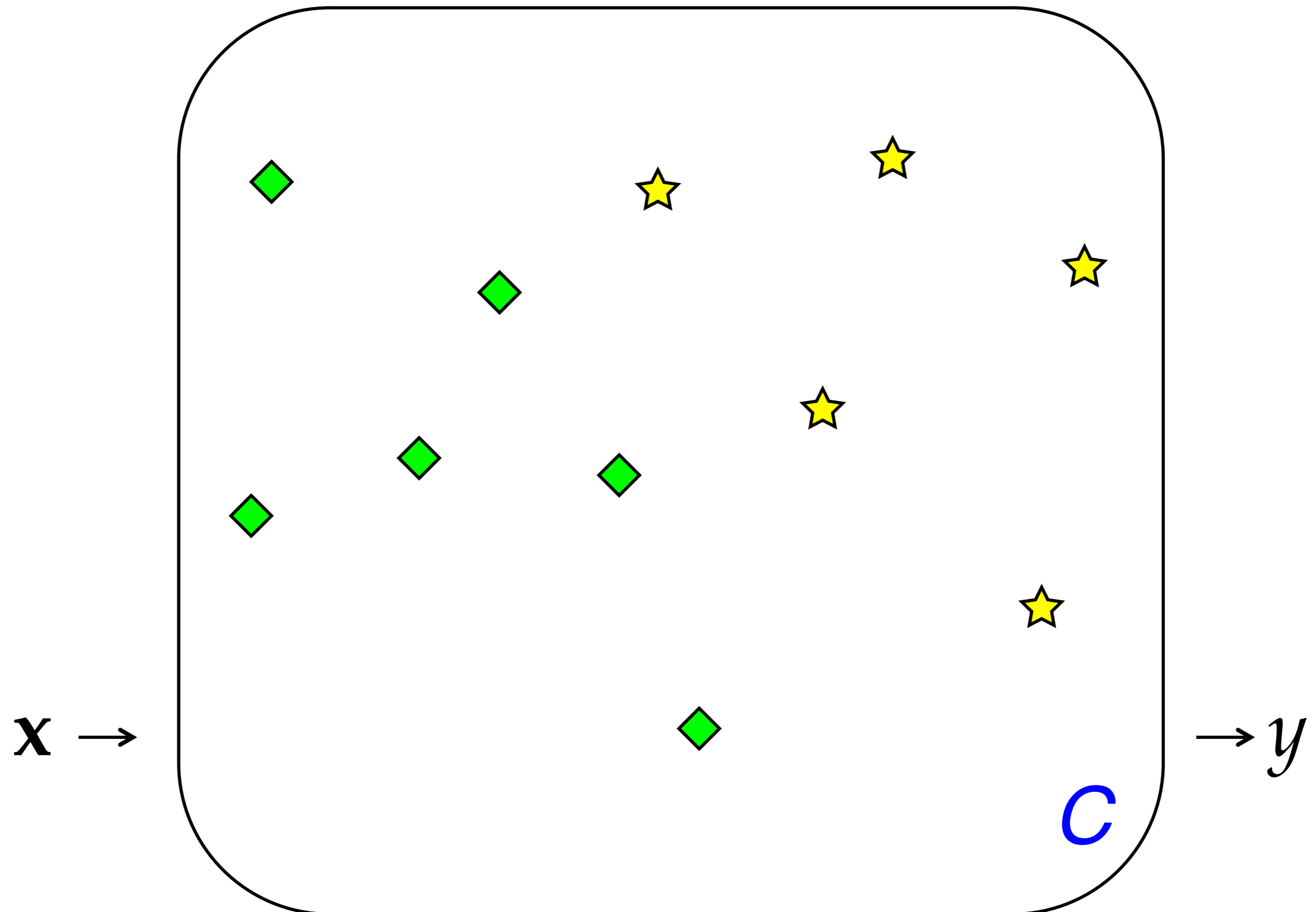
$\mathbf{x}'$  = Pet 's best **interest** in mind , but vets must follow law

- Weights are learned from data
- For the moment, assume binary classification: **financial** or **nonfinancial**
  - More positive weights more indicative of **financial**.
  - $\mathbf{w}^T \phi(\mathbf{x}) = 6.59$ ,  $\mathbf{w}^T \phi(\mathbf{x}') = -6.74$

# More than 2 classes

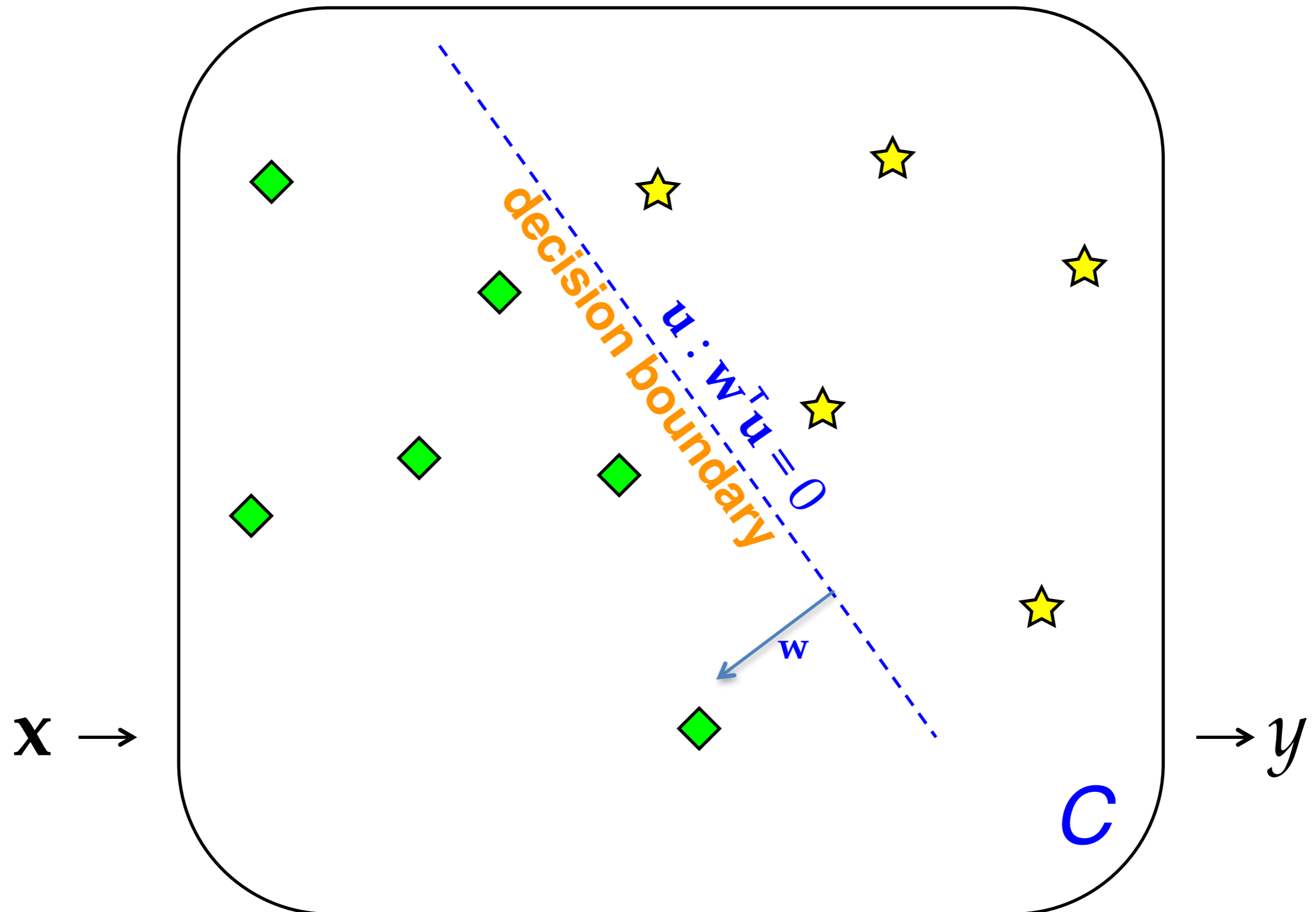
- Simply keep a separate weight vector for each label:  $\mathbf{w}_y$
- The label  $y$  whose weight vector gives the highest score wins!

# Linear Classifiers: Geometric View

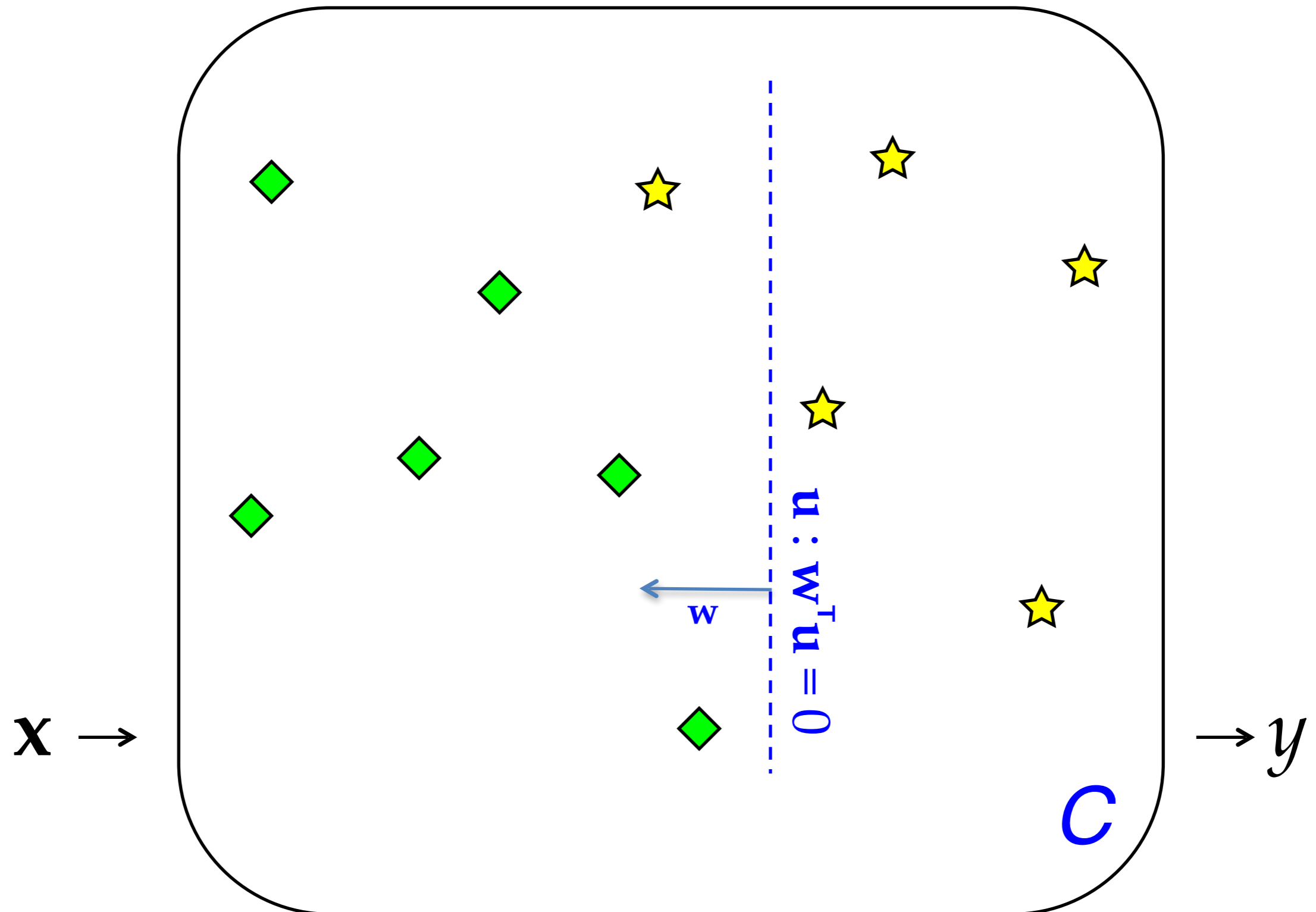




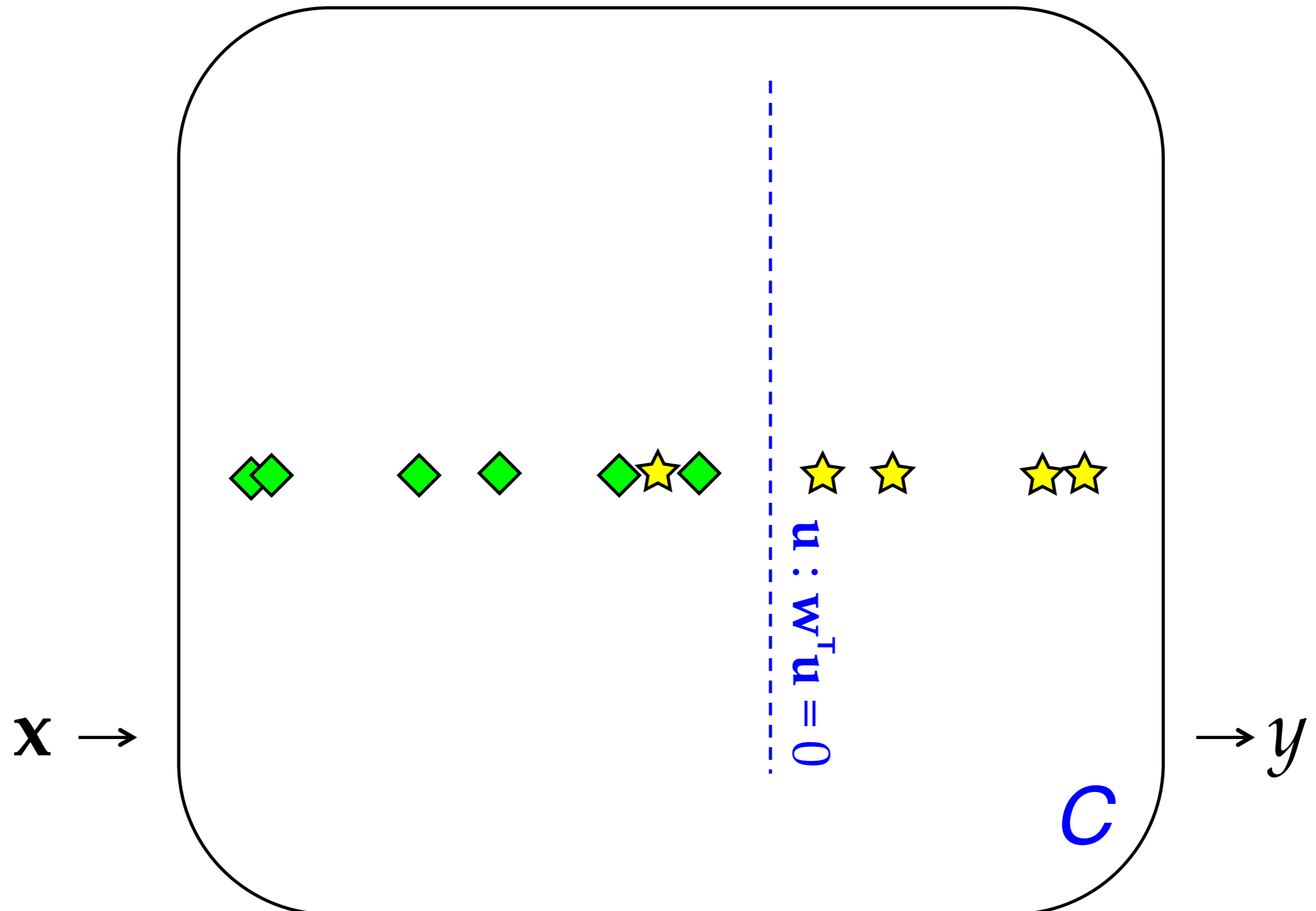
# Linear Classifiers: Geometric View



# Linear Classifiers: Geometric View



# Linear Classifiers: Geometric View



# Linear Classifiers (> 2 Classes)

return

$$\arg \max_y \mathbf{w}_y^\top \Phi(\mathbf{x})$$

$\mathbf{x} \rightarrow$

$\rightarrow y$

$C$

# The term “feature”

- The term “feature” is overloaded in NLP/ML. Here are three different concepts:
  - **Linguistic feature:** in some formalisms, a symbolic property that applies to a unit to categorize it, e.g. [–voice] for a sound in phonology or [+past] for a verb in morphology.
  - **Percept (or input feature):** captures some aspect of an input  $x$ ; binary- or real-valued. *[The term “percept” is nonstandard but I think it is useful!]*      **ends in -ing**
  - **Parameter (or model feature):** an association between some percept and an output class (or structure)  $y$  for which a real-valued weight or score is learned.      **ends in -ing  $\wedge$   $y$ =VERB**



# Weights

- So far we have just discussed the **classifier**, which relies on weights.
- The weights are determined by a **learner**, which fits them to the training data.
- Naïve Bayes probability estimation can be thought of as learning the weights with bag-of-words features only.
- Several popular learning algorithms for feature-based linear models: **perceptron**, support vector machine (SVM), maximum entropy a.k.a. multiclass logistic regression
  - Deep learning models are **nonlinear**

# Evaluating Multiclass Classifiers and Retrieval Algorithms

# Accuracy

- Assume we are disambiguating word senses such that every token has 1 gold sense label.
- The classifier predicts 1 label for each token in the test set.
- Thus, every test set token has a predicted label (*pred*) and a gold label (*gold*).
- The **accuracy** of our classifier is just the % of tokens for which the predicted label matched the gold label:  $\#_{pred=gold} / \#_{tokens}$



# Precision and Recall

- To measure the classifier with respect to a certain label  $y$ , and there are  $>2$ , we distinguish precision and recall:
  - **precision** = proportion of times the label was predicted and that prediction matched the gold:  $\#_{pred=gold=y} / \#_{pred=y}$
  - **recall** = proportion of times the label was in the gold standard and was recovered correctly by the classifier:  
 $\#_{pred=gold=y} / \#_{gold=y}$
- The harmonic mean of precision and recall, called **F<sub>1</sub>-score**, balances between the two.  
 $F_1 = 2 * precision * recall / (precision + recall)$

# Evaluating Retrieval Systems

- Precision/Recall/F-score are also useful for evaluating retrieval systems.
- E.g., consider a system which takes a word as input and is supposed to retrieve all rhymes.
- Now, for a single input (the query), there are often many correct outputs.
- **Precision** tells us whether most of the given outputs were valid rhymes; **recall** tells us whether most of the valid rhymes in the gold standard were recovered.

# Rhymes for “hinge”

**Gold**

**System**

klinge  
minge  
vinje

binge  
cringe  
fringe  
hinge  
impinge  
infringe  
syringe  
tinge  
twinge  
unhinge

ainge

# Rhymes for “hinge”

**Gold**

**System**

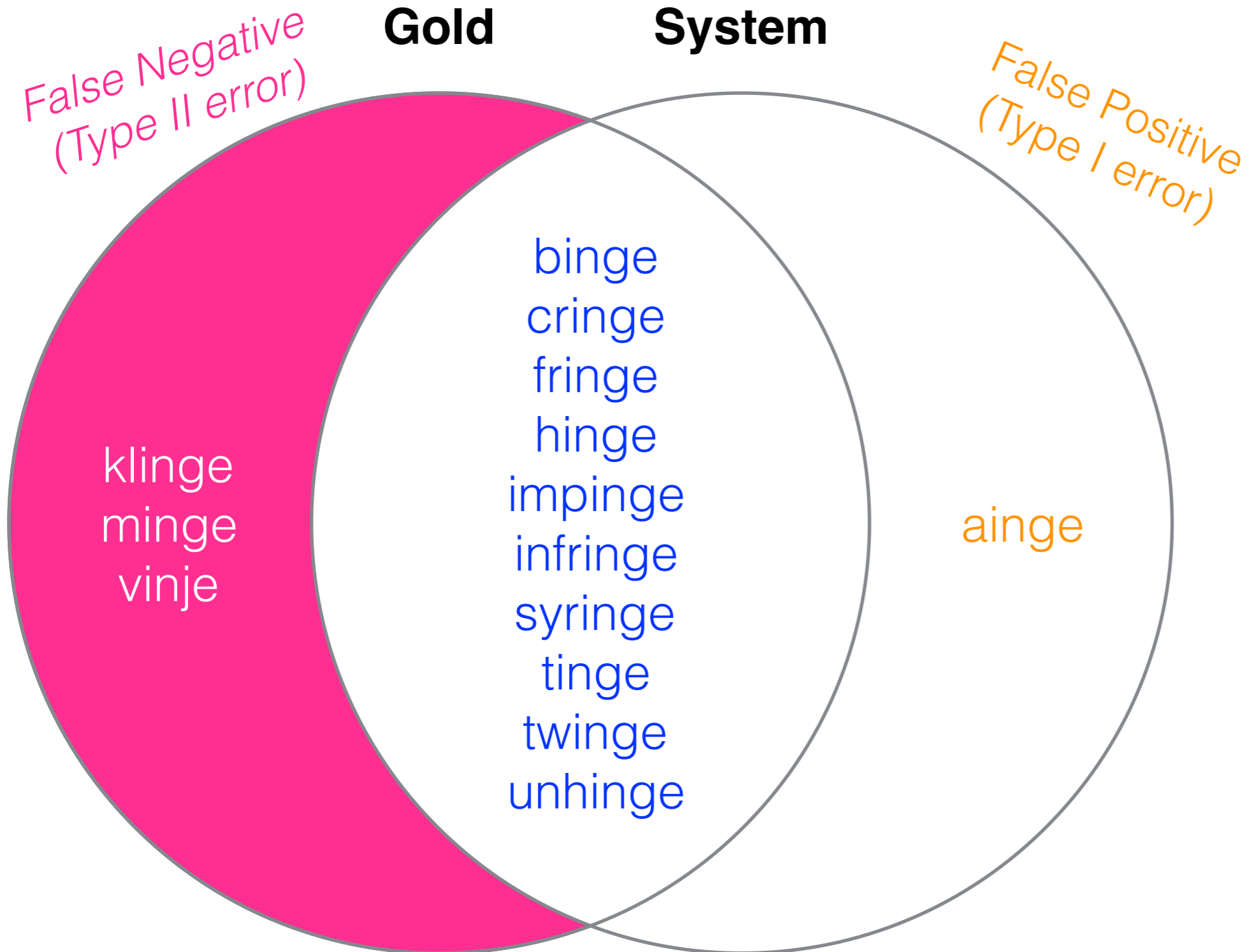
*False Positive  
(Type I error)*

klinge  
minge  
vinje

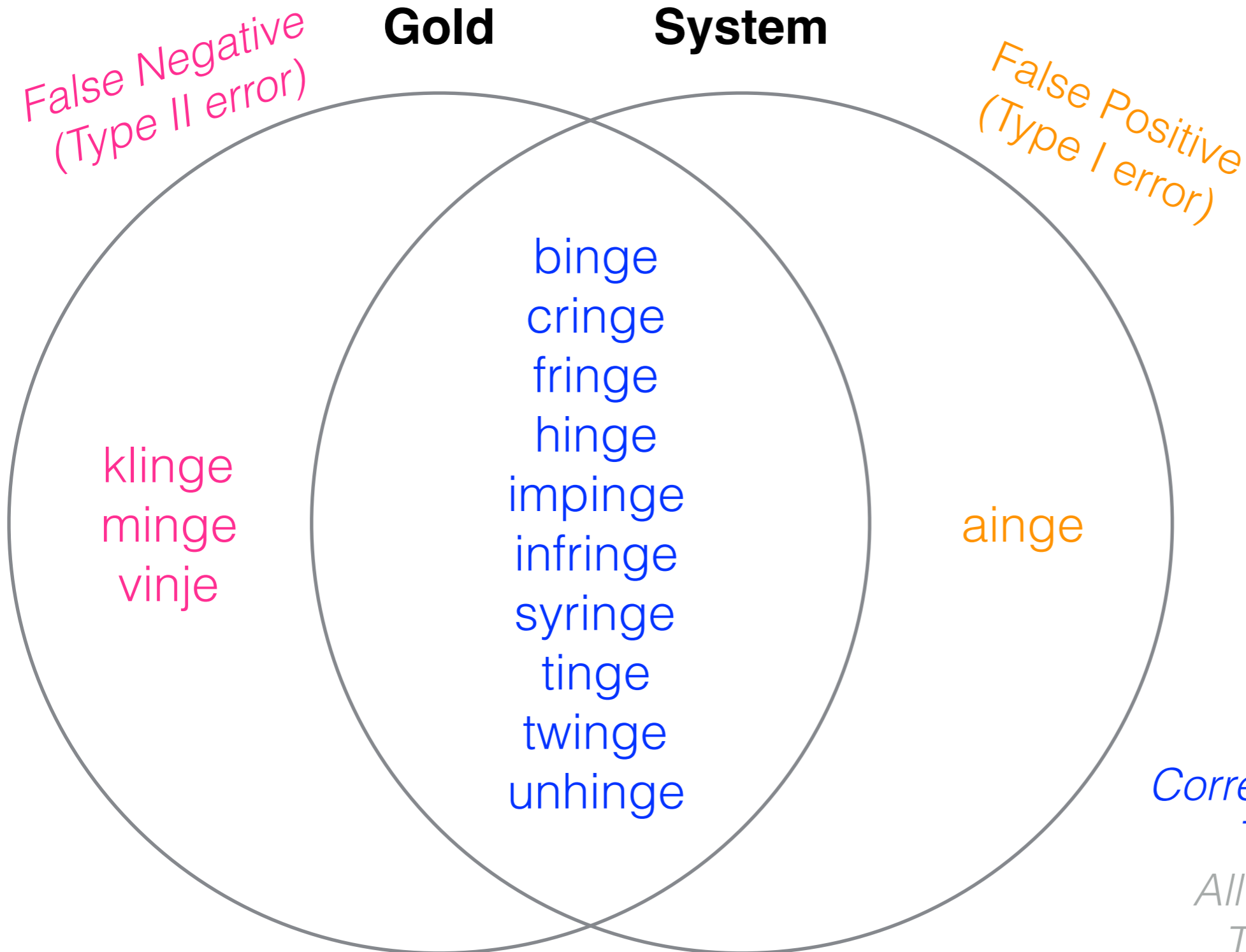
binge  
cringe  
fringe  
hinge  
impinge  
infringe  
syringe  
tinge  
twinge  
unhinge

ainge

# Rhymes for “hinge”



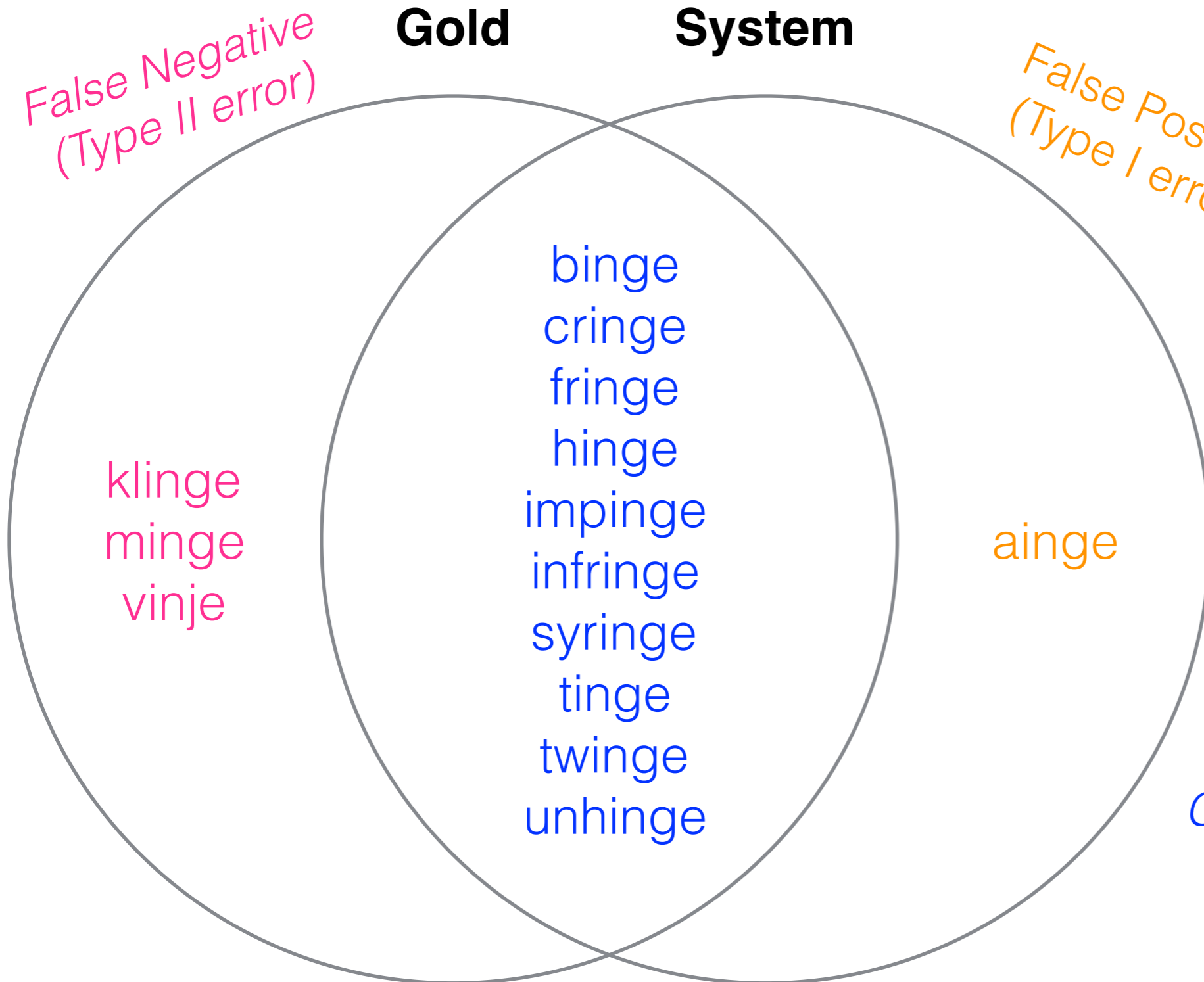
# Rhymes for “hinge”



	Sys=Y	Sys=N
Gold=Y	10	3
Gold=N	1	(large)

*Correctly predicted = True Positive*  
*All other words = True Negative*

# Precision & Recall



False Negative  
(Type II error)

False Positive  
(Type I error)

	Sys=Y	Sys=N
Gold=Y	10	3
Gold=N	1	(large)

**Precision = TP/(TP+FP)**  
 = 10/11 = 91%

**Recall = TP/(TP+FN)**  
 = 10/13 = 77%

**F<sub>1</sub> = 2·P·R/(P+R) = 83%**

*Correctly predicted = True Positive*  
*All other words = True Negative*

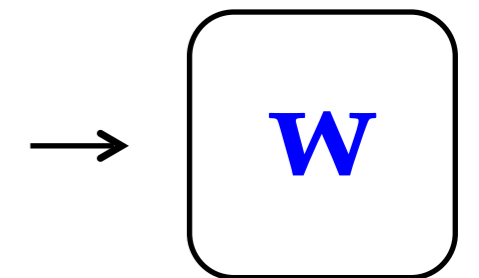
# Perceptron Learner



# Perceptron Learner

$\mathbf{X} \rightarrow$   
 $\mathbf{Y} \rightarrow$

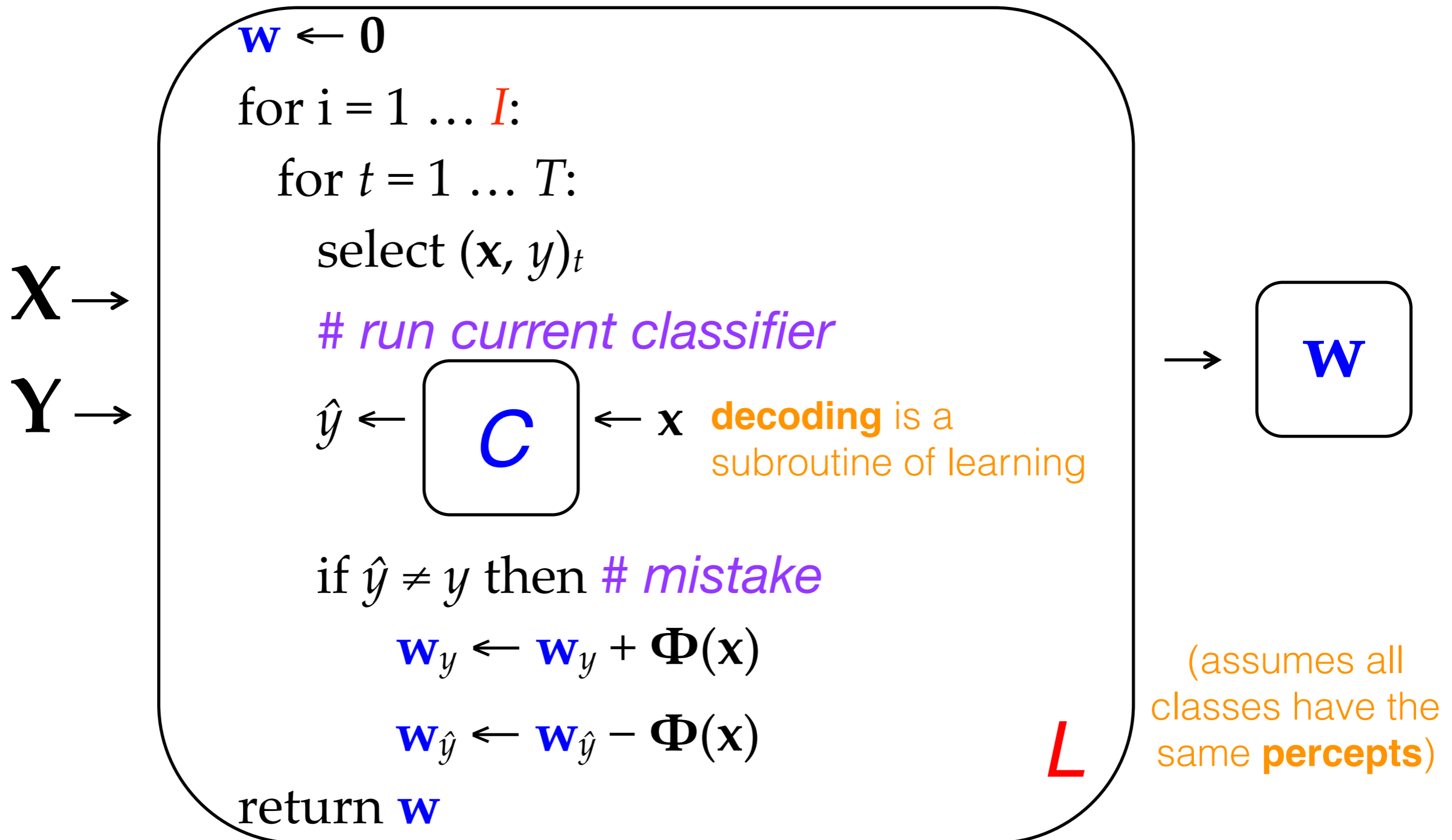
```
 $\mathbf{w} \leftarrow \mathbf{0}$   
for  $i = 1 \dots I$ :  
  for  $t = 1 \dots T$ :  
    select  $(\mathbf{x}, y)_t$   
    # run current classifier  
     $\hat{y} \leftarrow \arg \max_{y'} \mathbf{w}_{y'}^\top \Phi(\mathbf{x})$   
  
    if  $\hat{y} \neq y$  then # mistake  
       $\mathbf{w}_y \leftarrow \mathbf{w}_y + \Phi(\mathbf{x})$   
       $\mathbf{w}_{\hat{y}} \leftarrow \mathbf{w}_{\hat{y}} - \Phi(\mathbf{x})$   
  
return  $\mathbf{w}$ 
```



(assumes all classes have the same **percepts**)

**L**

# Perceptron Learner



work through example on the board

# Perceptron Learner

- **The perceptron doesn't estimate probabilities.** It just adjusts weights up or down until they classify the training data correctly.
  - No assumptions of feature independence necessary!  $\Rightarrow$  Better accuracy than NB
- The perceptron is an example of an **online** learning algorithm because it potentially updates its parameters (weights) with each training datapoint.
- Classification, a.k.a. **decoding**, is called with the latest weight vector. Mistakes lead to weight updates.
- One hyperparameter:  $I$ , the number of iterations (passes through the training data).
- Often desirable to make several passes over the training data. The number can be tuned. Under certain assumptions, it can be proven that the learner will converge.

# Perceptron: Avoiding overfitting

- Like any learning algorithm, the perceptron risks overfitting the training data. Two main techniques to improve generalization:
  - ▶ **Averaging:** Keep a copy of each weight vector as it changes, then average all of them to produce the final weight vector. [Daumé chapter](#) has a trick to make this efficient with large numbers of features.
  - ▶ **Early stopping:** Tune  $I$  by checking held-out accuracy on dev data (or cross-val on train data) after each iteration. If accuracy has ceased to improve, stop training and use the model from iteration  $I - 1$ .

# Generative vs. Discriminative

- Naïve Bayes allows us to classify via the **joint probability** of  $\mathbf{x}$  and  $y$ :
  - $p(y | \mathbf{x}) \propto p(y) \prod_{w \in \mathbf{x}} p(w | y)$   
=  $p(y) p(\mathbf{x} | y)$  (per the independence assumptions of the model)  
=  $p(y, \mathbf{x})$  (chain rule)
  - This means the model accounts for BOTH  $\mathbf{x}$  and  $y$ . From the joint distribution  $p(\mathbf{x}, y)$  it is possible to compute  $p(\mathbf{x})$  as well as  $p(y)$ ,  $p(\mathbf{x} | y)$ , and  $p(y | \mathbf{x})$ .
- NB is called a **generative** model because it assigns probability to linguistic objects ( $\mathbf{x}$ ). It could be used to generate “likely” language corresponding to some  $y$ . (Subject to its naïve modeling assumptions!)
  - (Not to be confused with the “generative” school of linguistics.)
- Some other linear models, including the perceptron, are **discriminative**: they are trained directly to classify given  $\mathbf{x}$ , and cannot be used to estimate the probability of  $\mathbf{x}$  or generate  $\mathbf{x} | y$ .

# Take-home points

- Feature-based linear classifiers are important to NLP.
  - You define the features, an algorithm chooses the weights.
    - \* The weights are real-valued.
    - \* Some classifiers, like **logistic regression**, are probabilistic: the weights correspond to probabilities.
  - More features  $\Rightarrow$  more flexibility, also more risk of overfitting. Because we work with large vocabularies, not uncommon to have millions of features.
- Some models, like Naïve Bayes, have a closed-form solution for parameters. Learning is cheap!
- The **perceptron** and other discriminative methods require fancier learning/optimization algorithms that iterate multiple times over the data, adjusting parameters until convergence (or some other stopping criterion).
  - The advantage: fewer modeling assumptions. Weights can be interdependent. **Discriminative** methods usually achieve higher accuracy with sufficient training data and computation (optimization).

# Which classifier to use?

- Fast and simple: **naïve Bayes**
- Very accurate, still simple: **perceptron**
- Very accurate, probabilistic, more complicated to implement: **MaxEnt**
- Potentially best accuracy, more complicated to implement: **SVM**
- All of these: watch out for **overfitting!** (NB—smoothing; Perceptron—early stopping, averaging; MaxEnt—regularization)
- Check the web for free and fast implementations, e.g. SVM<sup>light</sup>



# Further Reading: Basics & Examples

- **Manning:** features in linear classifiers  
<http://www.stanford.edu/class/cs224n/handouts/MaxentTutorial-16x9-FeatureClassifiers.pdf>
- **Goldwater:** naïve Bayes & MaxEnt examples  
[http://www.inf.ed.ac.uk/teaching/courses/fnlp/lectures/07\\_slides.pdf](http://www.inf.ed.ac.uk/teaching/courses/fnlp/lectures/07_slides.pdf)
- **O'Connor:** MaxEnt—incl. step-by-step examples, comparison to naïve Bayes  
<http://people.cs.umass.edu/~brenocon/inlp2015/04-logreg.pdf>
- **Daumé:** “The Perceptron” (*A Course in Machine Learning*, ch. 3)  
[http://www.ciml.info/dl/v0\\_8/ciml-v0\\_8-ch03.pdf](http://www.ciml.info/dl/v0_8/ciml-v0_8-ch03.pdf)
- **Neubig:** “The Perceptron Algorithm”  
<http://www.phontron.com/slides/nlp-programming-en-05-perceptron.pdf>

# Further Reading: Advanced

- **Neubig:** “Advanced Discriminative Learning”—MaxEnt w/ derivatives, SGD, SVMs, regularization  
<http://www.phontron.com/slides/nlp-programming-en-06-discriminative.pdf>
- **Manning:** generative vs. discriminative, MaxEnt likelihood function and derivatives  
<http://www.stanford.edu/class/cs224n/handouts/MaxentTutorial-16x9-MEMMs-Smoothing.pdf>, slides 3–20
- **Daumé:** linear models  
[http://www.ciml.info/dl/v0\\_8/ciml-v0\\_8-ch06.pdf](http://www.ciml.info/dl/v0_8/ciml-v0_8-ch06.pdf)
- **Smith:** A variety of loss functions for text classification  
<http://courses.cs.washington.edu/courses/cse517/16wi/slides/tc-intro-slides.pdf> & <http://courses.cs.washington.edu/courses/cse517/16wi/slides/tc-advanced-slides.pdf>