



LANGUAGE MODELS

Sandeep Soni

02/22/2024

CLASSIFICATION

x	y
The movie was fantastic	😊
The movie could have been better but wasn't bad	😐
That movie was a waste of time and money	😡

- Predict y from x
- $\hat{y} = \hat{h}(x)$
- Model \hat{h} as $P(Y|X)$
- X is high-dimensional
- $P(\text{😊} | \text{"the"}, \text{"movie"}, \text{"was"}, \text{"fantastic"})$

CLASSIFIERS

- Naive Bayes
 - Maximize $P(X, Y)$ and then use Bayes theorem
 - Conditional independence assumptions
$$P(X_i | Y) = P(X_i | Y, X_{i-1})$$
 - $P(X_i | Y)$ is estimated simply by counting
 - Smoothing is required for estimating probabilities

Generative

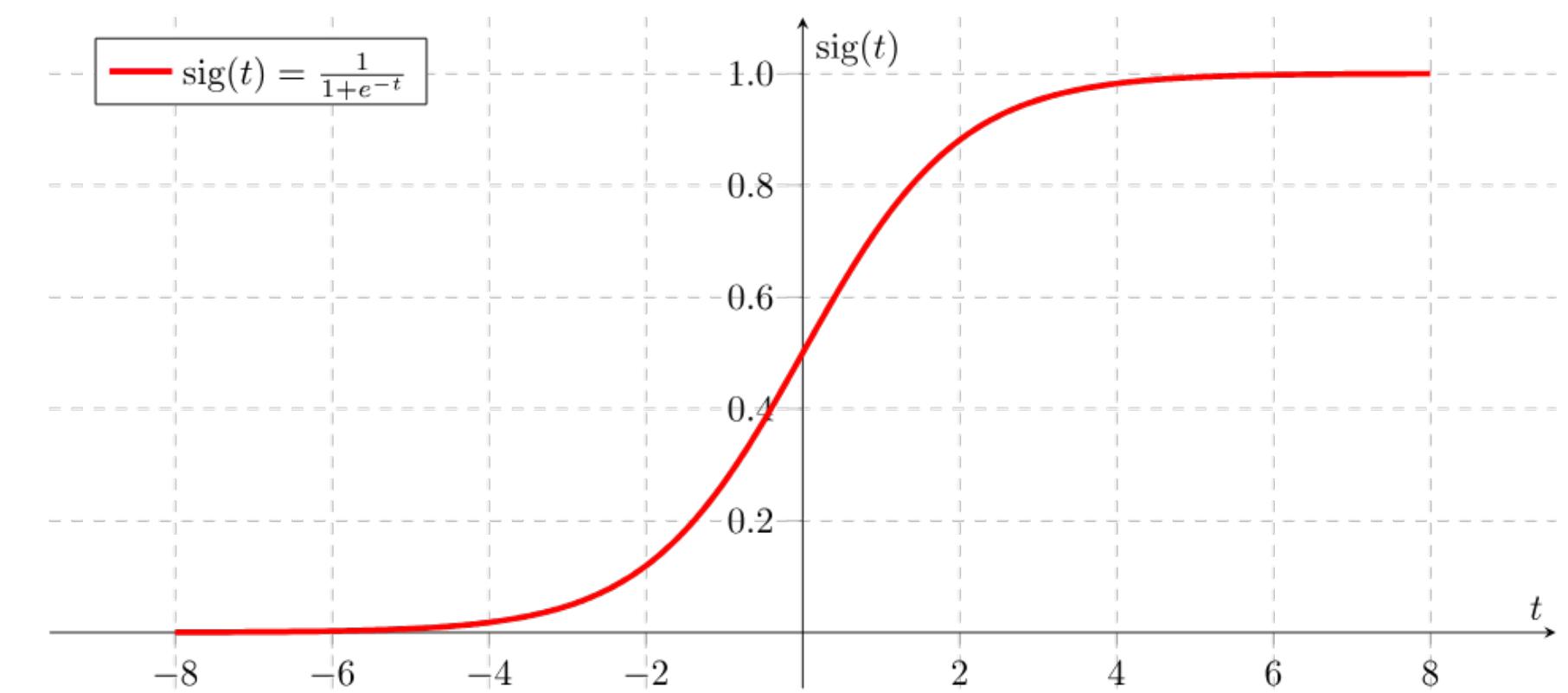
CLASSIFIERS

- Logistic Regression
 - Maximize $P(Y|X)$
 - $P(Y|X) = \sigma(\mathbf{w}^\top \mathbf{x})$
 - Weights learned using iterative algorithms
 - Regularization controlled with hyperparameters

Discriminative

CLASSIFIERS

- Logistic Regression
 - Maximize $P(Y|X)$
 - $P(Y|X) = \sigma(\mathbf{w}^\top \mathbf{x})$
 - Weights learned using iterative algorithms
 - Regularization controlled with hyperparameters



Discriminative

CLASSIFIERS

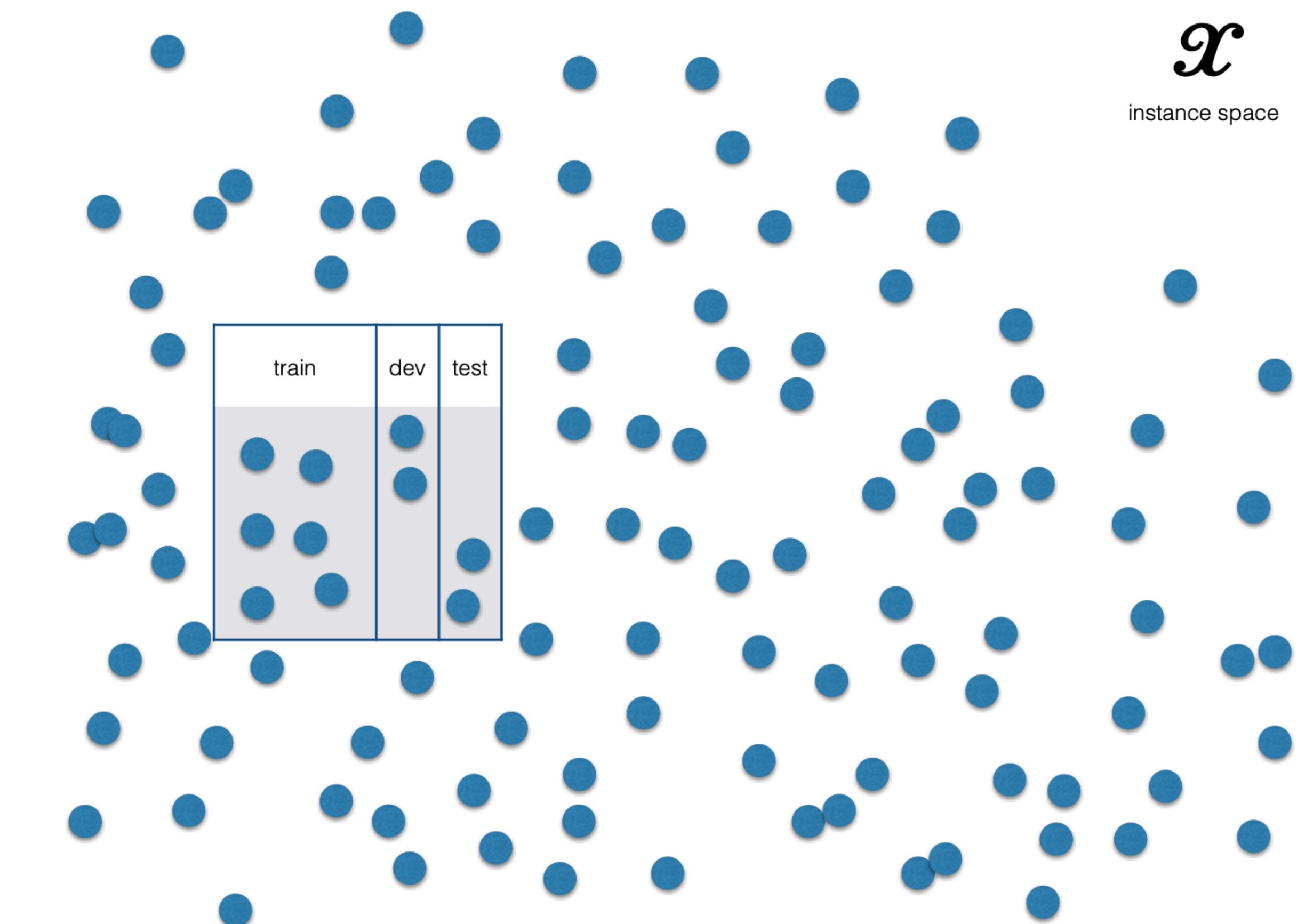
- Linear models: Naive Bayes, Logistic Regression
 - Perceptrons, support vector machines
- Non-linear models: Decision trees, kernel based methods, neural networks

SUPERVISED LEARNING

Train: Estimate weights, set hyperparameters

Dev: Evaluate using some performance metrics

Test: Do final evaluation



Slide from David Bamman

EVALUATION

	1	2	3	4	5	6	7	8	9	10
y	spam	spam	spam	spam	spam	ham	ham	ham	ham	ham
yhat	spam	spam	ham	spam	spam	spam	ham	spam	ham	spam

EVALUATION

	1	2	3	4	5	6	7	8	9	10
y	spam	spam	spam	spam	spam	ham	ham	ham	ham	ham
yhat	spam	spam	ham	spam	spam	spam	ham	spam	ham	spam

	y=spam	y=ham
yhat = spam	4	3
yhat=ham	1	2

EVALUATION

	1	2	3	4	5	6	7	8	9	10
y	spam	spam	spam	spam	spam	ham	ham	ham	ham	ham
yhat	spam	spam	ham	spam	spam	spam	ham	spam	ham	spam

	y=spam	y=ham
yhat = spam	True positives	False positives
yhat=ham	False negatives	True negatives

	y=spam	y=ham
yhat = spam	4	3
yhat=ham	1	2

EVALUATION

	1	2	3	4	5	6	7	8	9	10
y	spam	spam	spam	spam	spam	ham	ham	ham	ham	ham
yhat	spam	spam	ham	spam	spam	spam	ham	spam	ham	spam

	y=spam	y=ham
yhat = spam	True positives	False positives
yhat=ham	False negatives	True negatives

	y=spam	y=ham
yhat = spam	4	3
yhat=ham	1	2

$$N = \text{True positives} + \text{True negatives} + \text{False positives} + \text{False negatives}$$

EVALUATION METRICS

$$F_1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

EVALUATION METRICS

Predicted \ Observed		y=spam	y=ham
yhat = spam	True positives	False positives	
yhat=ham	False negatives	True negatives	

$$F_1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

EVALUATION METRICS

$$\text{Accuracy} = \frac{\text{True positives} + \text{True negatives}}{N}$$

Predicted \ Observed		y=spam	y=ham
yhat = spam	True positives	False positives	
yhat=ham	False negatives	True negatives	

$$F_1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

EVALUATION METRICS

$$\text{Accuracy} = \frac{\text{True positives} + \text{True negatives}}{N}$$

$$\text{Precision} = \frac{\text{True positives}}{\text{True positives} + \text{False positives}}$$

$$F_1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

		y=spam	y=ham
Predicted \ Observed	yhat = spam	True positives	False positives
	yhat=ham	False negatives	True negatives

EVALUATION METRICS

$$\text{Accuracy} = \frac{\text{True positives} + \text{True negatives}}{N}$$

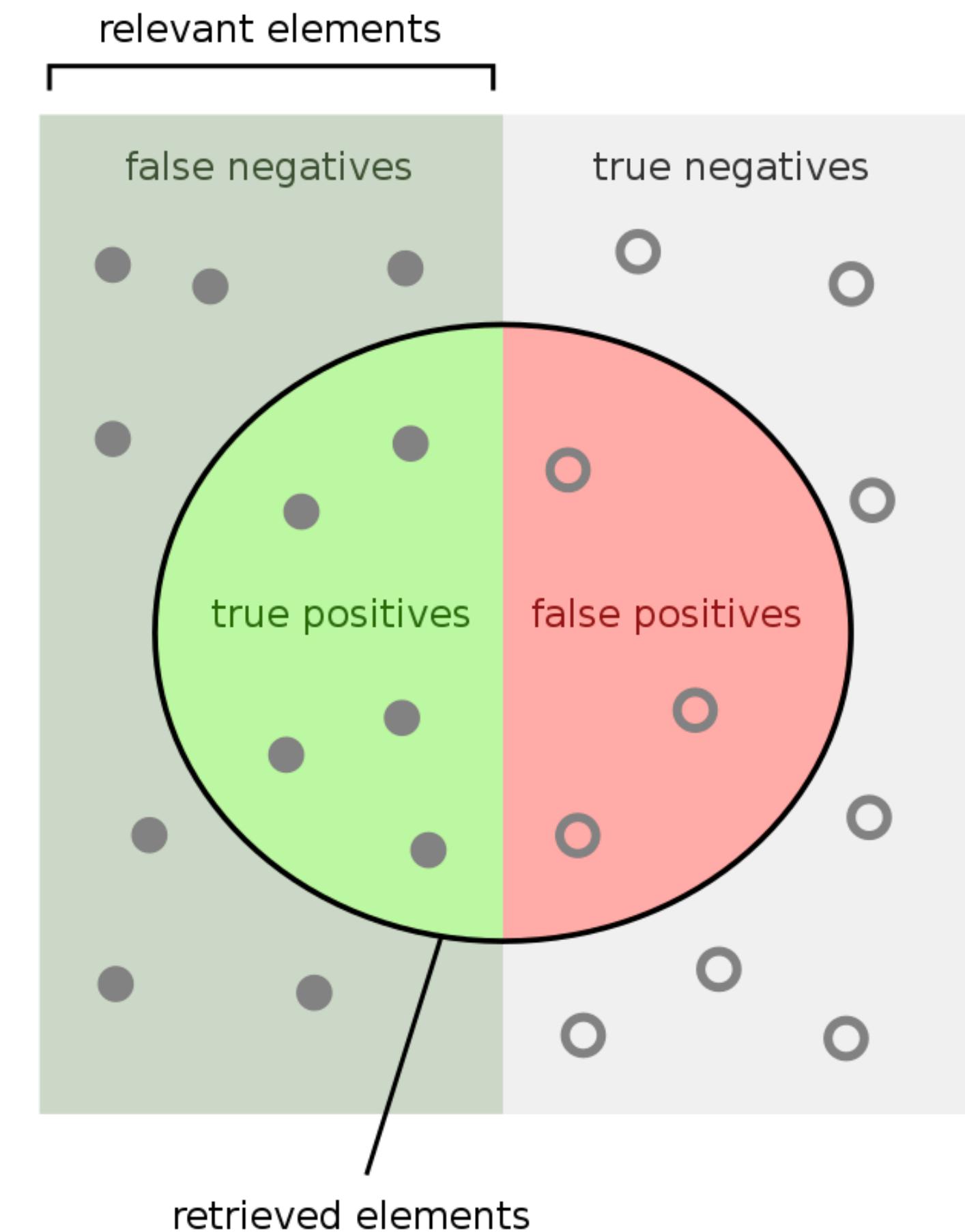
$$\text{Precision} = \frac{\text{True positives}}{\text{True positives} + \text{False positives}}$$

$$\text{Recall} = \frac{\text{True positives}}{\text{True positives} + \text{False positives}}$$

$$F_1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

		Predicted \ Observed	y=spam	y=ham
yhat = spam	yhat = spam	True positives	False positives	
	yhat=ham	False negatives	True negatives	

PRECISION AND RECALL



$$\text{Precision} = \frac{\text{How many retrieved items are relevant?}}{\text{How many relevant items are retrieved?}}$$
$$\text{Recall} = \frac{\text{How many retrieved items are relevant?}}{\text{How many relevant items are there?}}$$

Source: Wikipedia

IN CLASS

- Text classification demo

CLASSIFICATION

$$P(y|x) = \frac{P(y)P(x|y)}{P(x)}$$

Bayes Theorem

$$P(y|x) = \frac{P(y)P(x|y)}{\sum_{y'} P(x,y')}$$

$$P(y|x) = \frac{P(y)P(x|y)}{\sum_{y'} P(y')P(x|y')}$$

Marginalization

What is the probability of a given piece of text i.e. $P(x)$?

WHY SHOULD WE MODEL $P(x)$?

- NLP perspective: many tasks are sequence-to-sequence tasks
 - Machine translation: text in source language \rightarrow text in target language
 - Summarization: long text \rightarrow short text

WHY SHOULD WE MODEL $P(x)$?

- A system should be “fluent” in generating output sequences
- How can we quantify fluency?

El cafe negro me gusta mucho

Machine
translation
system I

The coffee black me pleases much

El cafe negro me gusta mucho

Machine
translation
system II

I love dark coffee

Machine
translation
system I

The coffee black me pleases much

$P(\text{"The coffee black me pleases much"})$

<

Machine
translation
system II

I love dark coffee

$P(\text{"I love dark coffee"})$

“When I look at an article in Russian, I say: `This is really written in English but it has been coded in some strange symbols. I will now proceed to decode.”

–Warren Weaver (1955)

NOISY CHANNEL MODEL

$$P_{\text{english} \leftarrow \text{spanish}}(\mathbf{w}^e | \mathbf{w}^s) \propto P_{\text{english}, \text{spanish}}(\mathbf{w}^e, \mathbf{w}^s)$$

$$= P_{\text{spanish} \leftarrow \text{english}}(\mathbf{w}^{(s)} | \mathbf{w}^{(e)}) \times P_{\text{english}}(\mathbf{w}^{(e)})$$

↓ ↓

Translation model Language model

We want the noisy channel model to have high fidelity to the source text and high fluency in the generated output

LANGUAGE MODEL

- \mathcal{V} is a vocabulary of symbols that need to be generated
- \mathcal{S} is a infinite set of sequences of symbols taken from \mathcal{V} ; each sequence ends with a special token **STOP**
- A sequence $x \in \mathcal{S}$

LANGUAGE MODEL

- $P(x) = P(x_1, x_2, \dots, x_n)$
- $P(\text{"I will be back"}) = P(\text{"I"}, \text{"will"}, \text{"be"}, \text{"back"}, \text{STOP})$
- $0 \leq P(x) \leq 1$
- $\sum_{x \in \mathcal{S}} P(x') = 1$

LANGUAGE MODEL

- Language modeling is a task of estimating a probabilistic model over words i.e. $P(x)$
- This is not easy!
- What is the $P(\text{"to be or not to be"})$

CHAIN RULE

$$P(x_1, x_2, \dots, x_n) = P(x_1)$$

$$P(\text{to be or not to be}) = P(\text{to})$$

$$\times P(x_2 | x_1)$$

$$\times P(\text{be} | \text{to})$$

$$\times P(x_3 | x_1, x_2)$$

$$\times P(\text{or} | \text{to, be})$$

...

...

$$\times P(x_n | x_1, \dots, x_{n-1})$$

$$\times P(\text{be} | \text{to, be, ..., to})$$

CHAIN RULE

$$P(\text{to be or not to be}) = P(\text{to})$$

$$\times P(\text{be}|\text{to})$$

$$\times P(\text{or}|\text{to, be})$$

...

$$\times P(\text{be}|\text{to, be, ..., to})$$

$P(\text{"be"} \mid \text{"to be or not to"})$
is hard to estimate unless
we make some assumptions

MARKOV ASSUMPTIONS

$$P(x_i | x_{i-1}, \dots, x_1) \approx P(x_i)$$

zero-order

$$P(x_i | x_{i-1}, \dots, x_1) \approx P(x_i | x_{i-1})$$

first-order

$$P(x_i | x_{i-1}, \dots, x_1) \approx P(x_i | x_{i-2}, x_{i-1})$$

second-order

N-GRAM LANGUAGE MODELS

$$P(x) = \prod_i P(x_i) \times P(\text{STOP})$$

unigram

$$P(x) = \prod_i P(x_i | x_{i-1}) \times P(\text{STOP} | x_n)$$

bigram

$$P(x) = \prod_i P(x_i | x_{i-2}, x_{i-1}) \times P(\text{STOP} | x_{n-1}, x_n)$$

trigram

BIGRAM GENERATION

$$P(x) = \prod_i P(x_i | x_{i-1}) \times P(\text{STOP} | x_n)$$

bigram

Previous word	Current word	Probability
START	I	0.2
START	like	0.0008
...
I	like	0.1
I	coffee	0.0003
...
coffee	STOP	0.06
I	STOP	0.0001

START _____

BIGRAM GENERATION

$$P(x) = \prod_i P(x_i | x_{i-1}) \times P(\text{STOP} | x_n)$$

bigram

Previous word	Current word	Probability
START	I	0.2
START	like	0.0008
...
I	like	0.1
I	coffee	0.0003
...
coffee	STOP	0.06
I	STOP	0.0001

START I

START I _____

BIGRAM GENERATION

$$P(x) = \prod_i P(x_i | x_{i-1}) \times P(\text{STOP} | x_n)$$

bigram

Previous word	Current word	Probability
START	I	0.2
START	like	0.0008
...
I	like	0.1
I	coffee	0.0003
...
coffee	STOP	0.06
I	STOP	0.0001

START I

START I like

START I like _____

BIGRAM GENERATION

$$P(x) = \prod_i P(x_i | x_{i-1}) \times P(\text{STOP} | x_n)$$

bigram

Previous word	Current word	Probability
START	I	0.2
START	like	0.0008
...
I	like	0.1
I	coffee	0.0003
...
coffee	STOP	0.06
I	STOP	0.0001

START I

START I like

START I like black

START I like black _____

BIGRAM GENERATION

$$P(x) = \prod_i P(x_i | x_{i-1}) \times P(\text{STOP} | x_n)$$

bigram

Previous word	Current word	Probability
START	I	0.2
START	like	0.0008
...
I	like	0.1
I	coffee	0.0003
...
coffee	STOP	0.06
I	STOP	0.0001

START I

START I like

START I like black

START I like black coffee

START I like black coffee _____

BIGRAM GENERATION

$$P(x) = \prod_i P(x_i | x_{i-1}) \times P(\text{STOP} | x_n)$$

bigram

Previous word	Current word	Probability
START	I	0.2
START	like	0.0008
...
I	like	0.1
I	coffee	0.0003
...
coffee	STOP	0.06
I	STOP	0.0001

START I

START I like

START I like black

START I like black coffee

START I like black coffee STOP

ESTIMATION

$$P(x) = \prod_i P(x_i) \times P(\text{STOP})$$

unigram

$$P(x) = \prod_i P(x_i | x_{i-1}) \times P(\text{STOP} | x_n)$$

bigram

$$P(x) = \prod_i P(x_i | x_{i-2}, x_{i-1}) \times P(\text{STOP} | x_{n-1}, x_n)$$

trigram

Maximum likelihood
estimation

$$P(x_i) = \frac{c(x_i)}{N}$$

$$P(x_i | x_{i-1}) = \frac{c(x_i, x_{i-1})}{c(x_{i-1})}$$

$$P(x_i | x_{i-2}, x_{i-1}) = \frac{c(x_i, x_{i-1}, x_{i-2})}{c(x_{i-1}, x_{i-2})}$$

DATA SPARSITY

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

Figure 4.1 Bigram counts for eight of the words (out of $V = 1446$) in the Berkeley Restaurant Project corpus of 9332 sentences. Zero counts are in gray.

SLP 4.3

If probabilities are estimated from a fixed finite sample, we'll encounter situations where we won't see some terms, bigrams, etc?

DATA SPARSITY

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

Figure 4.1 Bigram counts for eight of the words (out of $V = 1446$) in the Berkeley Restaurant Project corpus of 9332 sentences. Zero counts are in gray.

SLP 4.3

How to calculate $P(x)$? How to estimate $P(x_i)$?

SMOOTHING

Maximum likelihood estimate

$$P(x_i) = \frac{c(x_i)}{N}$$

Smoothed estimate

$$P(x_i) = \frac{c(x_i) + \alpha}{N + V\alpha}$$

Various different ways to do smoothing; other solutions include interpolation

EVALUATION

- How do we know if our language model is good?
- Ideally, we should evaluate a language model extrinsically (e.g., how much a language model improves the performance of a machine translation model)
- More often we evaluate a language model intrinsically

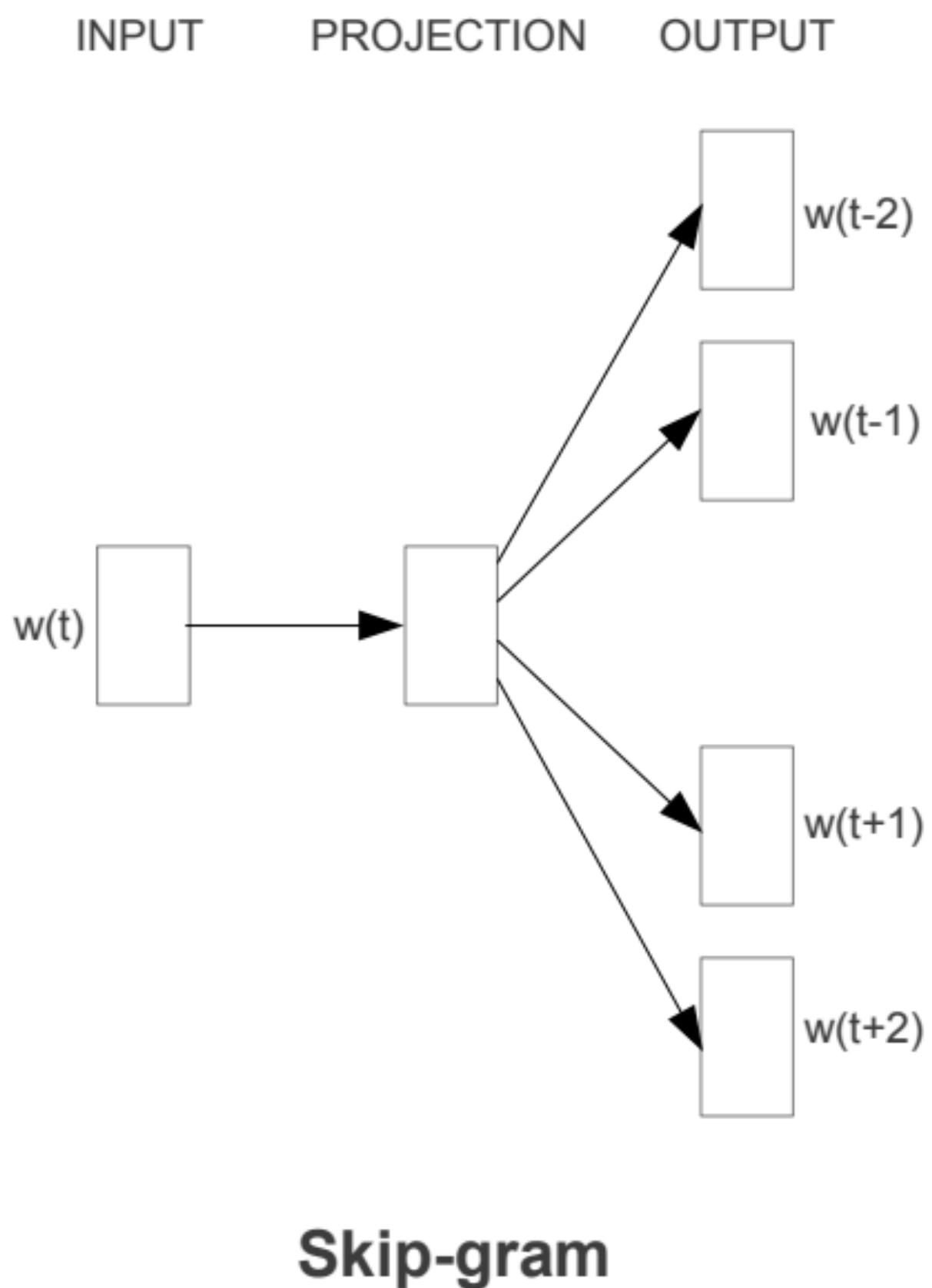
EVALUATION

- A good language model should assign high probability to unseen sequences that are plausible
- A good language model should be less perplexed or surprised
- Perplexity = inverse probability of test data, averaged over words

PERPLEXITY

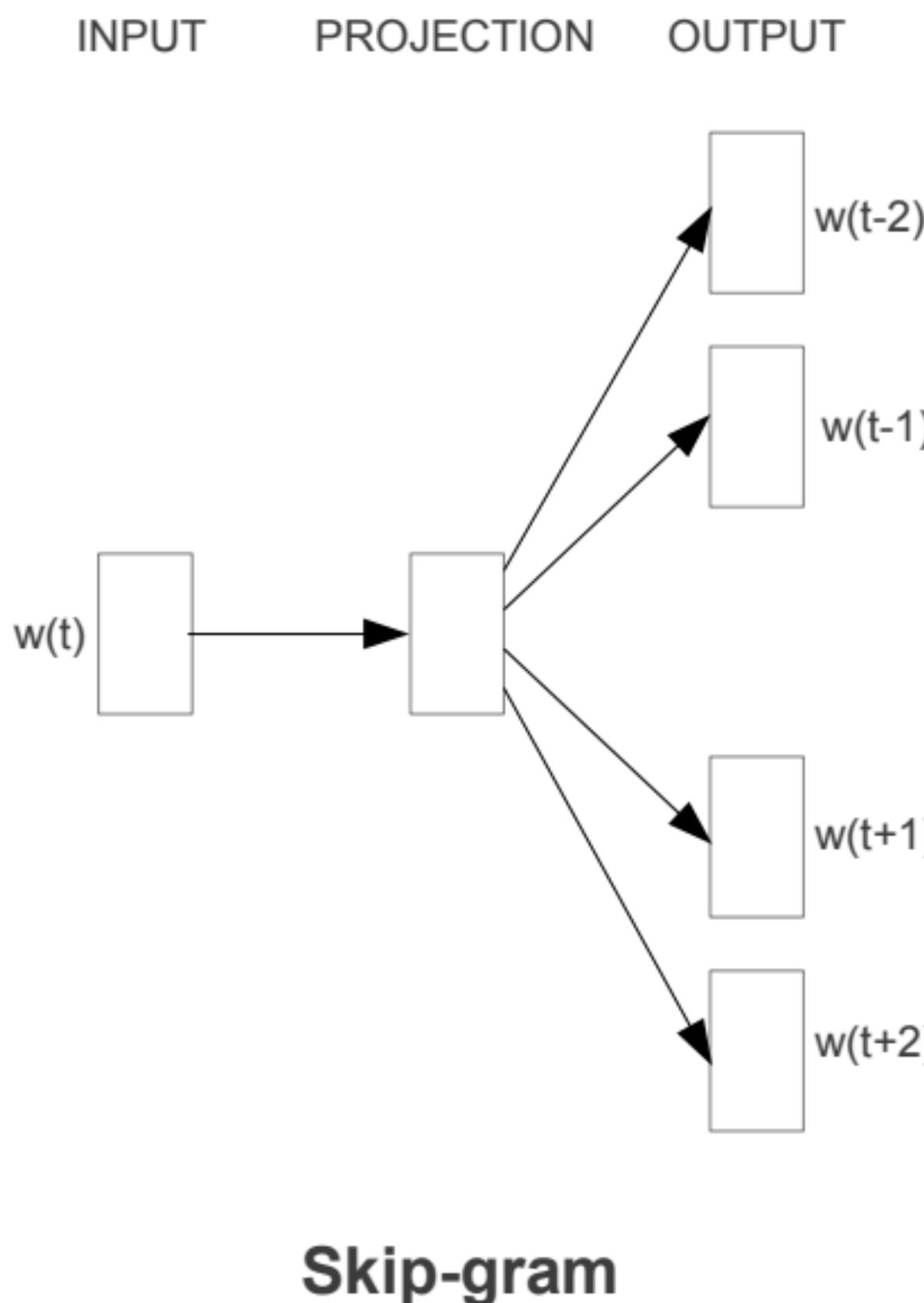
- Perplexity = $2^{-\frac{l(w)}{M}}$, where M is the total number of unseen tokens
- $$l(w) = \sum_{m=1}^M \log P(w_m | w_{m-1}, \dots, w_1)$$
- Based on the language model of your choice, you'll calculate $l(w)$
- Smaller perplexity is better!

SKIPGRAM



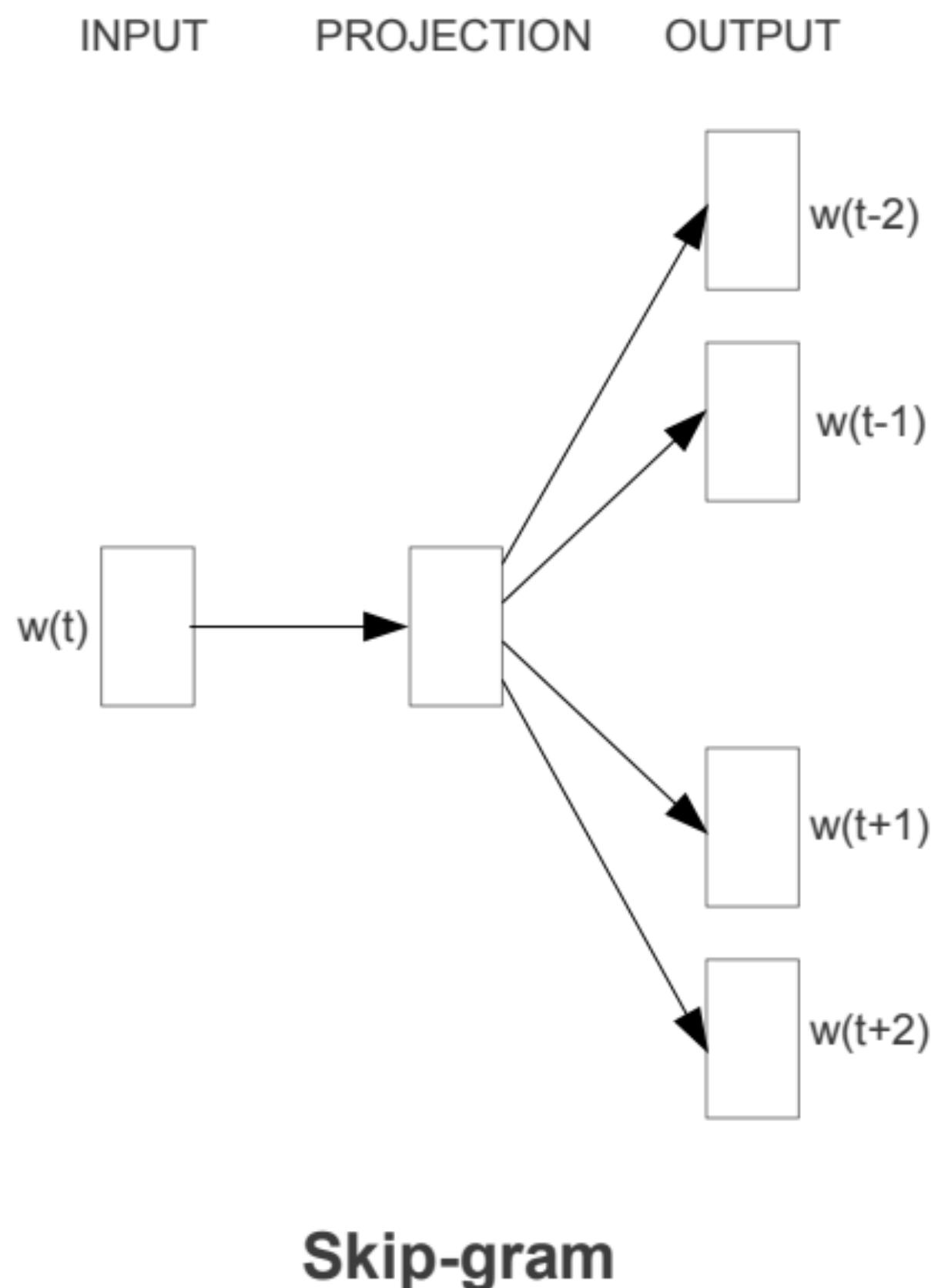
- In one of the most sensational comebacks in the champions league Liverpool defeated Barcelona by mounting a **strong** performance in the second half

SKIPGRAM



- In one of the most sensational comebacks in the champions league Liverpool defeated Barcelona by mounting a **strong** _____ in the second half

SKIPGRAM



- In one of the most sensational comebacks in the champions league Liverpool defeated Barcelona by mounting a **strong** _____ in the second half

Word2vec learned word vectors using the skipgram language model

IN CLASS

- lm exploration