

POLYNOMIAL AND NON LINEAR REGRESSION

Ariana Villegas

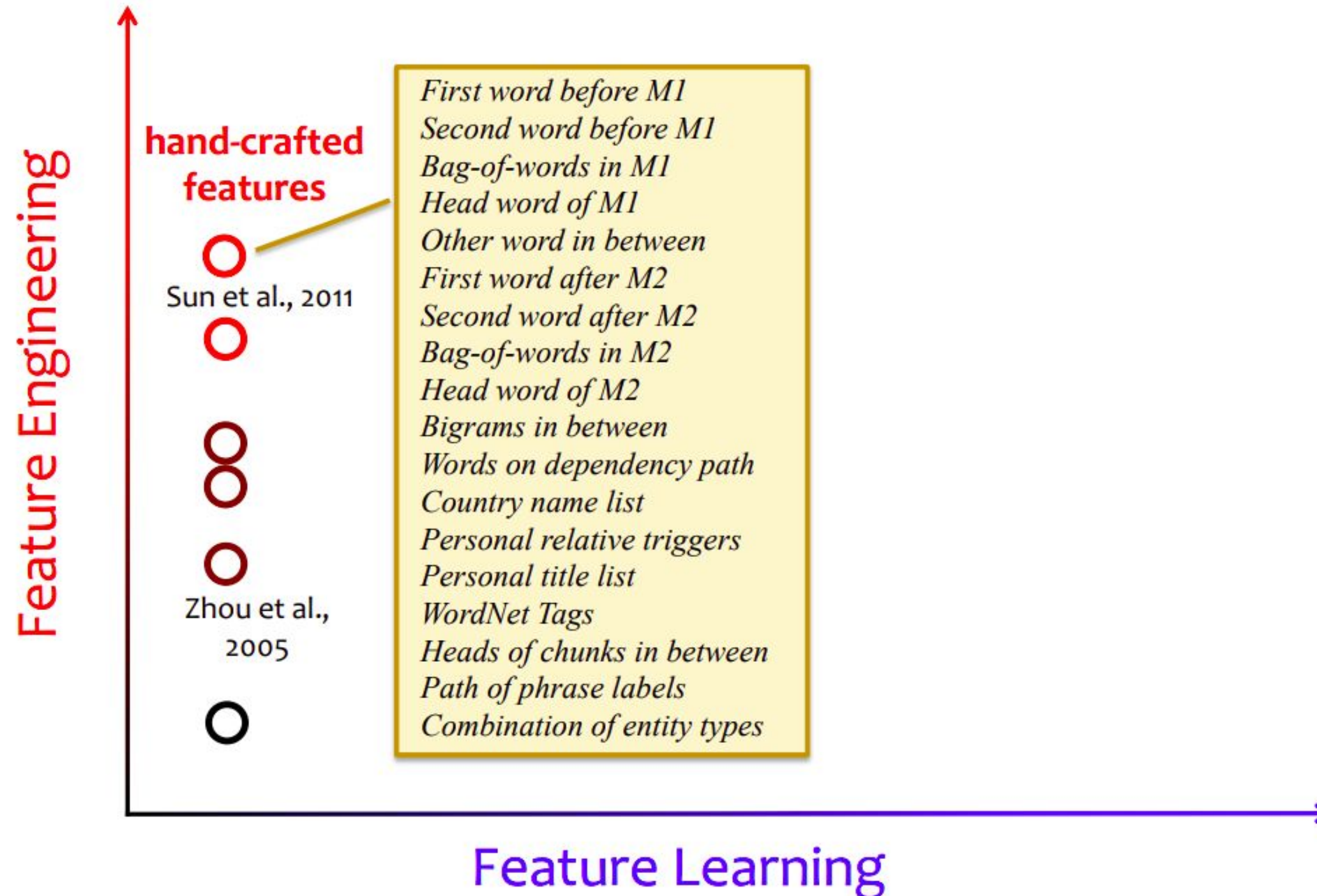
1.



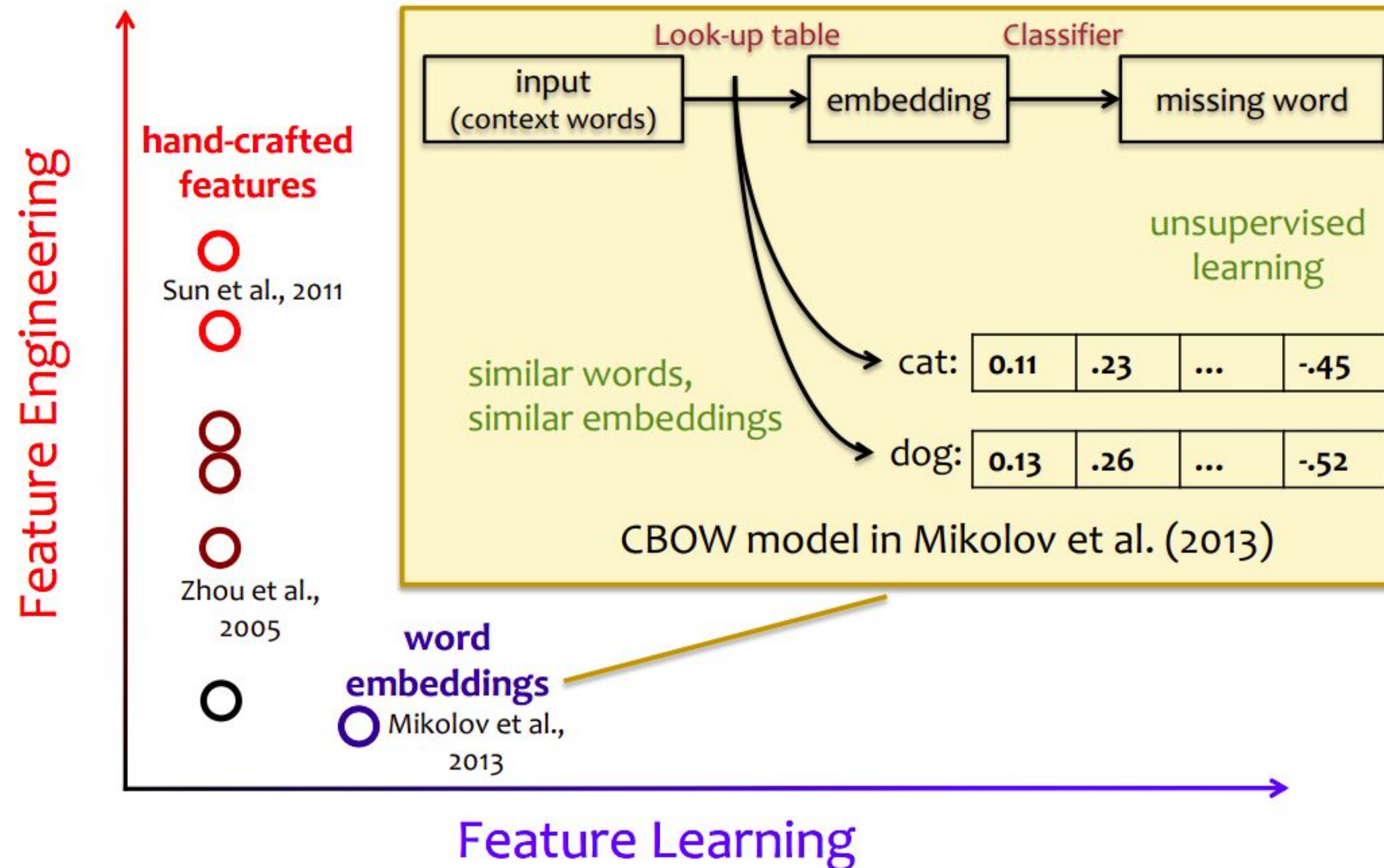
Feature Engineering

NLP & CV

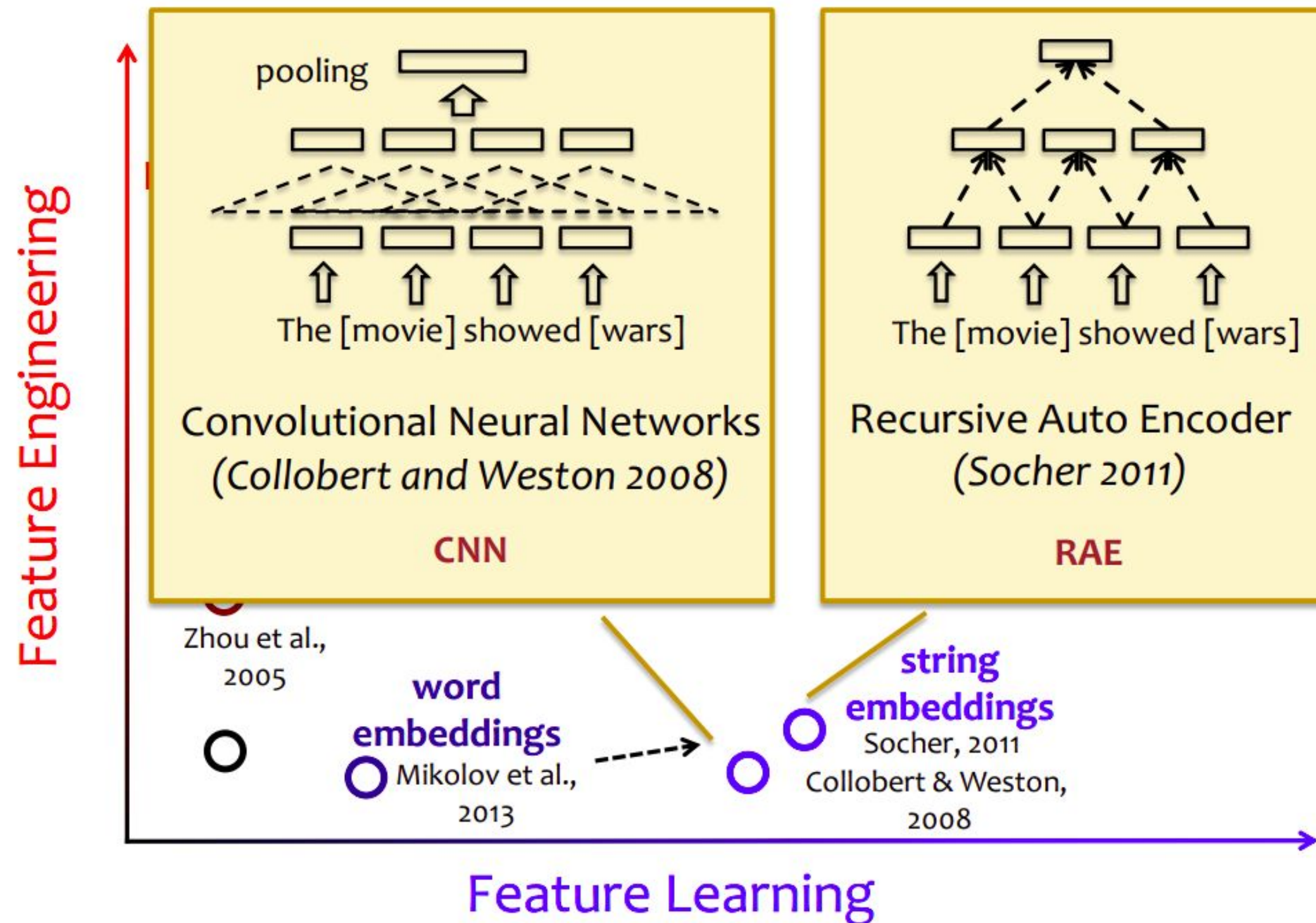
Where do features come from?



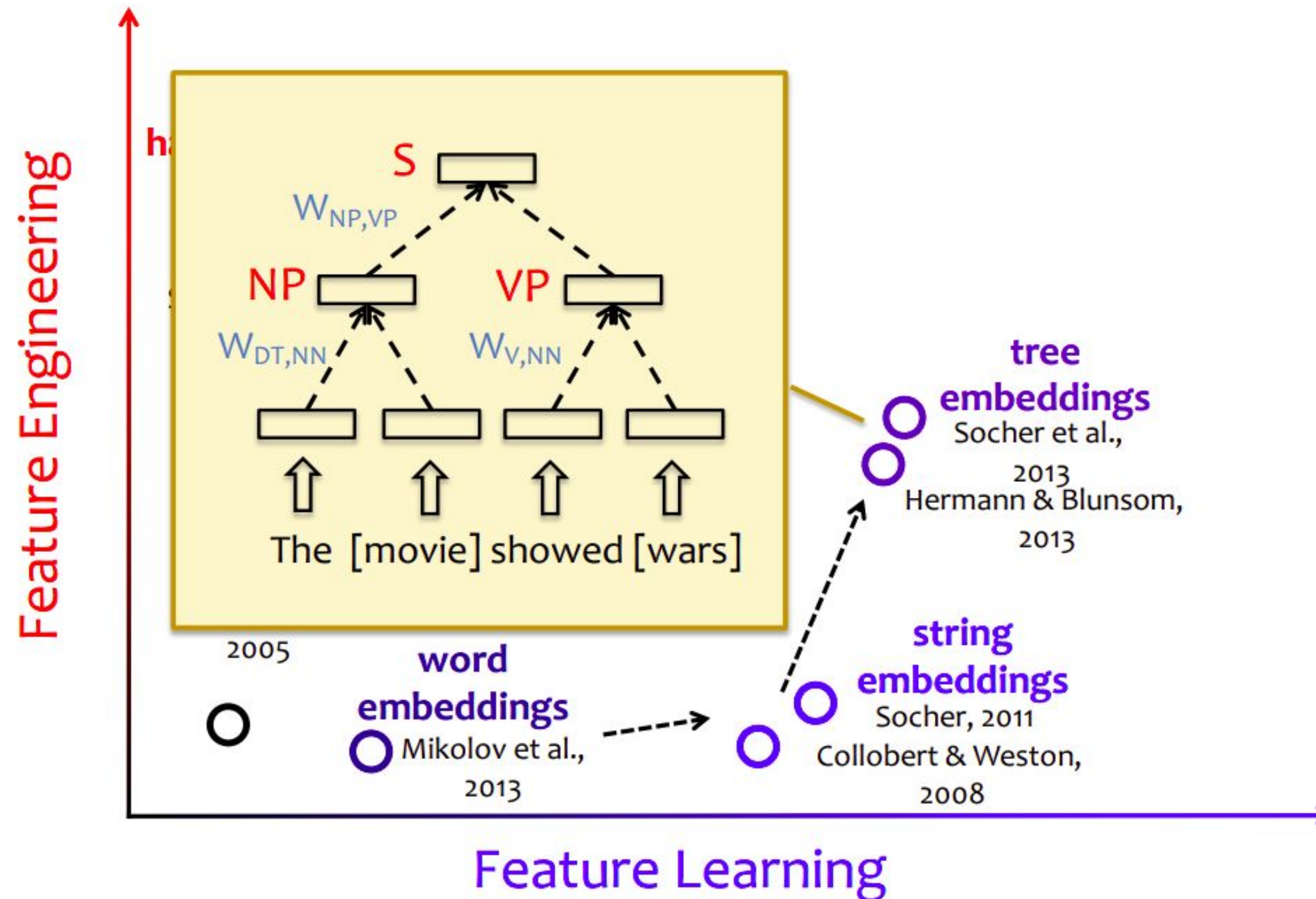
Where do features come from?



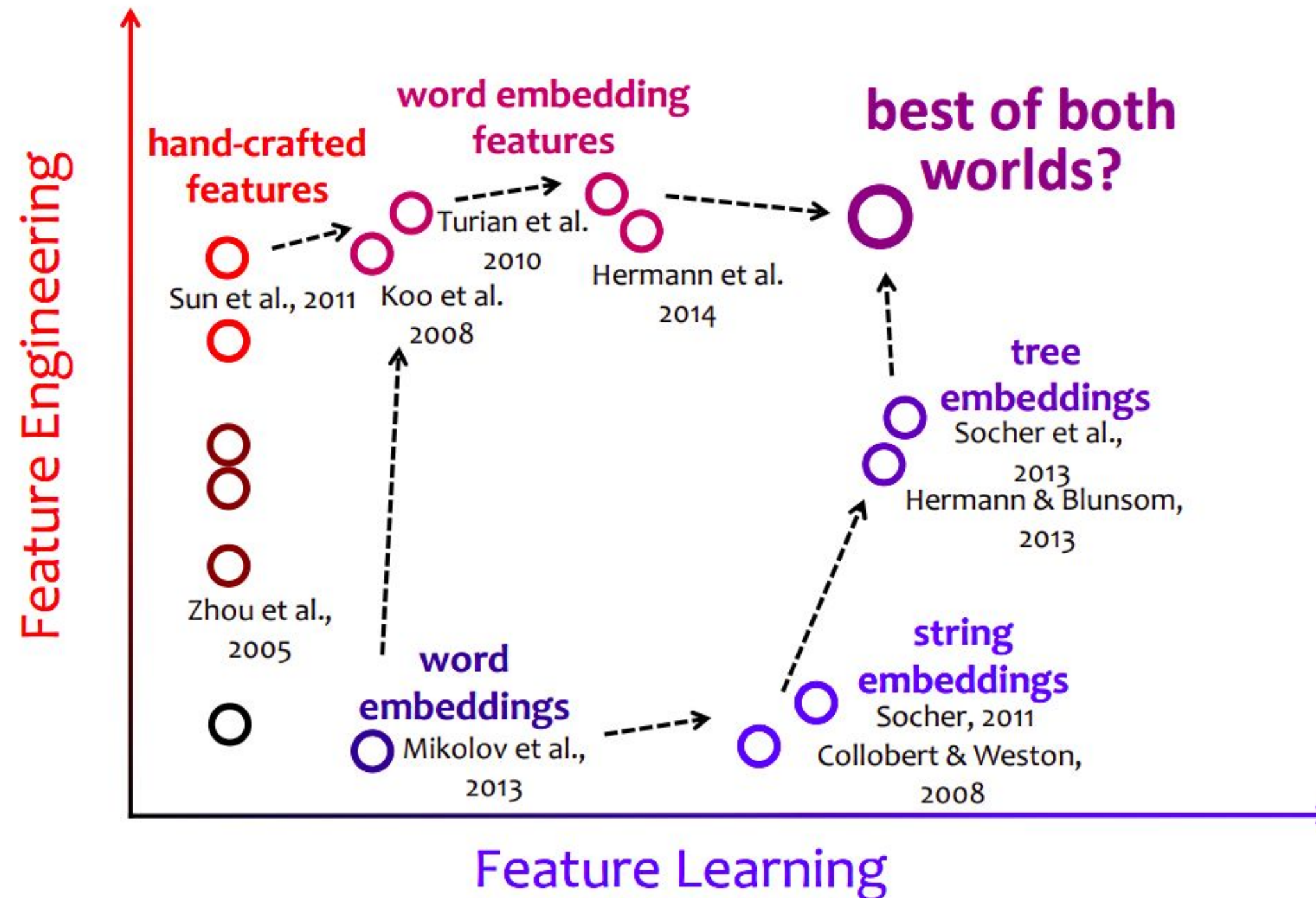
Where do features come from?



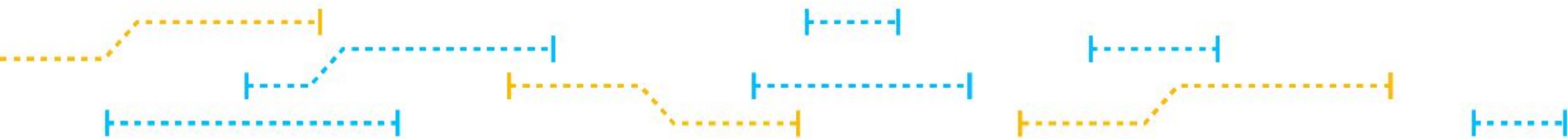
Where do features come from?



Where do features come from?



What features should you use for text?



What features should you use for text?

Per-word Features:

	$x^{(1)}$	$x^{(2)}$	$x^{(3)}$	$x^{(4)}$	$x^{(5)}$	$x^{(6)}$
<code>is-capital(w_i)</code>	1	0	1	0	0	0
<code>endswith(w_i, "e")</code>	1	1	0	0	0	1
<code>endswith(w_i, "d")</code>	0	0	0	1	1	0
<code>endswith(w_i, "ed")</code>	0	0	0	1	1	0
<code>w_i == "aardvark"</code>	0	0	0	0	0	0
<code>w_i == "hope"</code>	0	0	0	0	0	1
...

deter. noun noun verb verb noun
The movie I watched depicted hope



What features should you use for text?

Context Features:

	$x^{(1)}$	$x^{(2)}$	$x^{(3)}$	$x^{(4)}$	$x^{(5)}$	$x^{(6)}$
...
$w_i == \text{"watched"}$	0	0	0	1	0	0
$w_{i+1} == \text{"watched"}$	0	0	1	0	0	0
$w_{i-1} == \text{"watched"}$	0	0	0	0	1	0
$w_{i+2} == \text{"watched"}$	0	1	0	0	0	0
$w_{i-2} == \text{"watched"}$	0	0	0	0	0	1
...

Table 3. Tagging accuracies with different feature templates and other changes on the *WSJ* 19-21 development set.

Model	Feature Templates	# Feats	Sent. Acc.	Token Acc.	Unk. Acc.
3GRAMMEMM	See text	248,798	52.07%	96.92%	88.99%
NAACL 2003	See text and [1]	460,552	55.31%	97.15%	88.61%
Replication	See text and [1]	460,551	55.62%	97.18%	88.92%
Replication'	+rareFeatureThresh = 5	482,364	55.67%	97.19%	88.96%
5W	+ $\langle t_0, w_{-2} \rangle, \langle t_0, w_2 \rangle$	730,178	56.23%	97.20%	89.03%
5WSHAPES	+ $\langle t_0, s_{-1} \rangle, \langle t_0, s_0 \rangle, \langle t_0, s_{+1} \rangle$	731,661	56.52%	97.25%	89.81%
5WSHAPESDS	+ distributional similarity	737,955	56.79%	97.28%	90.46%

deter. noun noun verb verb noun

The movie I watched depicted hope

What features should you use for text?

Table 3. Tagging accuracies with different feature templates and other changes on the *WSJ* 19-21 development set.

Model	Feature Templates	# Feats	Sent. Acc.	Token Acc.	Unk. Acc.
3GRAMMEMM	See text	248,798	52.07%	96.92%	88.99%
NAACL 2003	See text and [1]	460,552	55.31%	97.15%	88.61%
Replication	See text and [1]	460,551	55.62%	97.18%	88.92%
Replication'	+rareFeatureThresh = 5	482,364	55.67%	97.19%	88.96%
5W	+ $\langle t_0, w_{-2} \rangle, \langle t_0, w_2 \rangle$	730,178	56.23%	97.20%	89.03%
5WSHAPES	+ $\langle t_0, s_{-1} \rangle, \langle t_0, s_0 \rangle, \langle t_0, s_{+1} \rangle$	731,661	56.52%	97.25%	89.81%
5WSHAPESDS	+ distributional similarity	737,955	56.79%	97.28%	90.46%

deter. noun noun verb verb noun
The movie I watched depicted hope



Word Embeddings

One-hot vectors

- Standard representation of a word in NLP: 1-hot vector (aka. a string)
- Vectors representing related words share nothing in common

	a	and	be	cat	dog	...	you	zebra
cat:	0	0	0	1	0	...	0	0
dog:	0	0	0	0	1	...	0	0

Word embeddings

- Word embedding: real-valued vector representation of a word in M dimensions
- Related words have similar vectors
- Long history in NLP: Term-doc frequency matrices, Reduce dimensionality with {LSA, NNMF, CCA, PCA}, Brown clusters, Vector space models, Random projections, Neural networks / deep learning

cat:	0.13	.26	...	-.52
dog:	0.11	.23	...	-.45

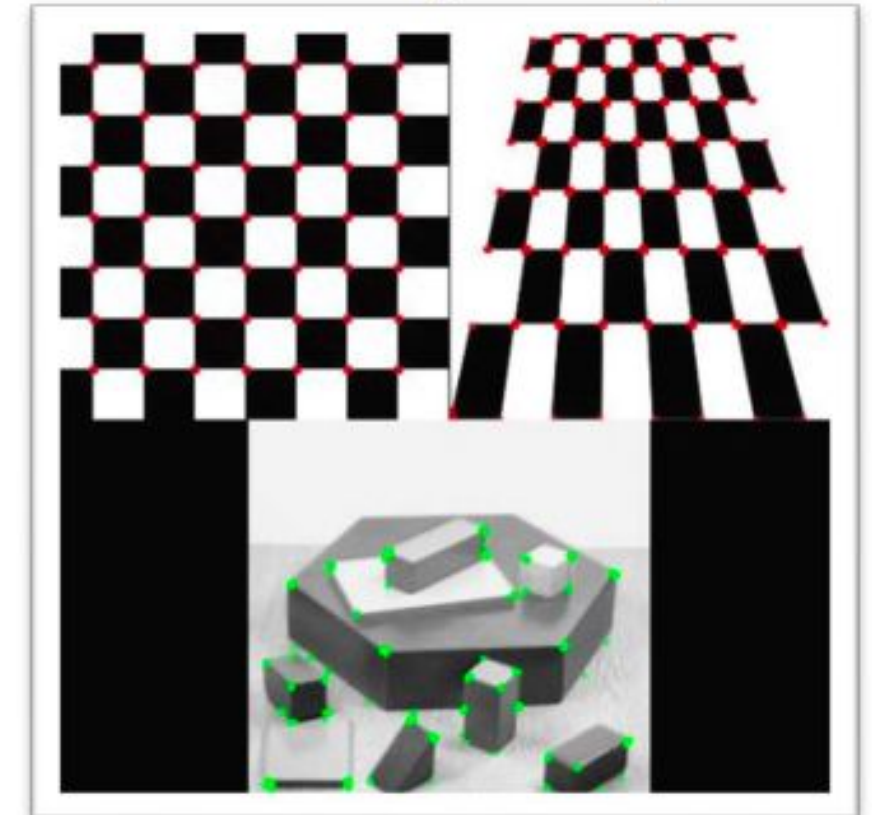


What features should you use for images?

Edge detection (Canny)



Corner Detection (Harris)



2.



Non-linear Regression

Definition & Examples

**What if our linear model is not
good? How can we create a
more
complicated model?**

Fitting a Polynomial

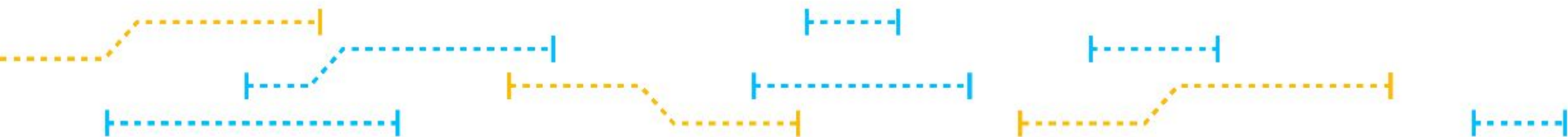
We can create a more complicated model by defining input variables that are **combinations of components of x**

Example: an **M-th order polynomial** function of **one dimensional** feature **x**:

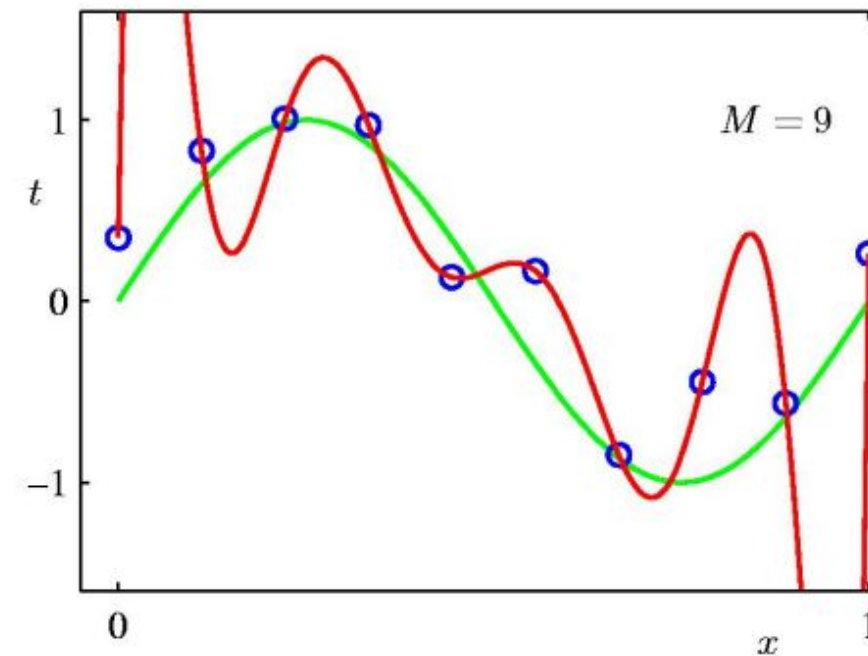
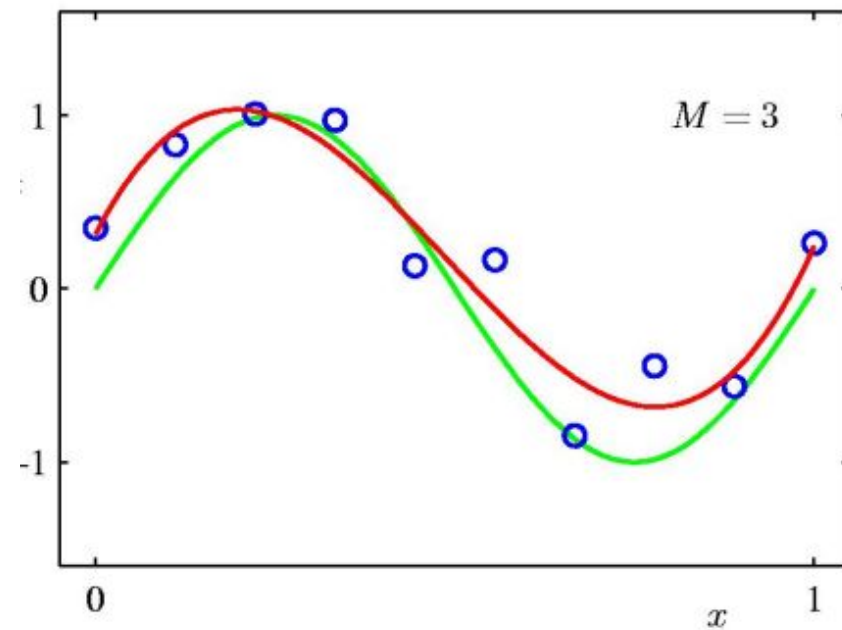
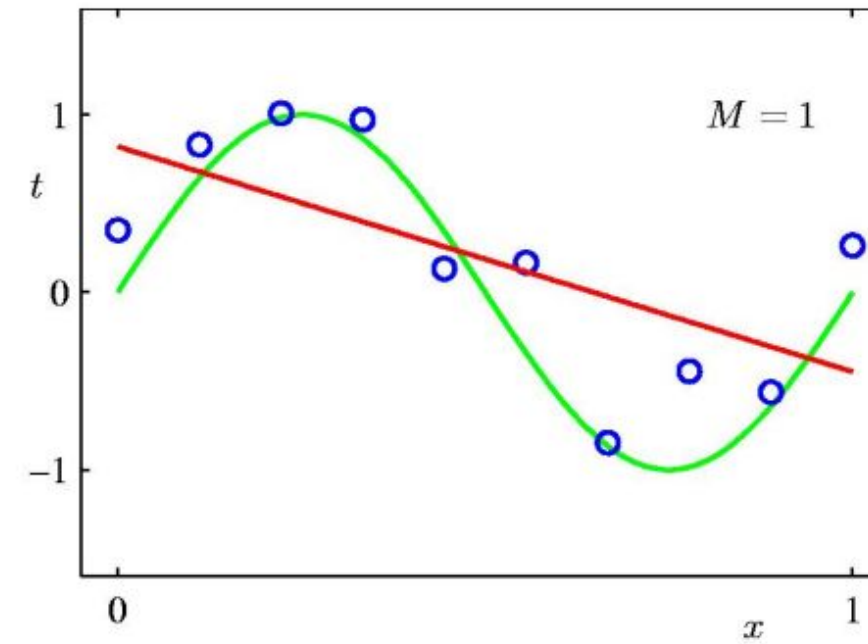
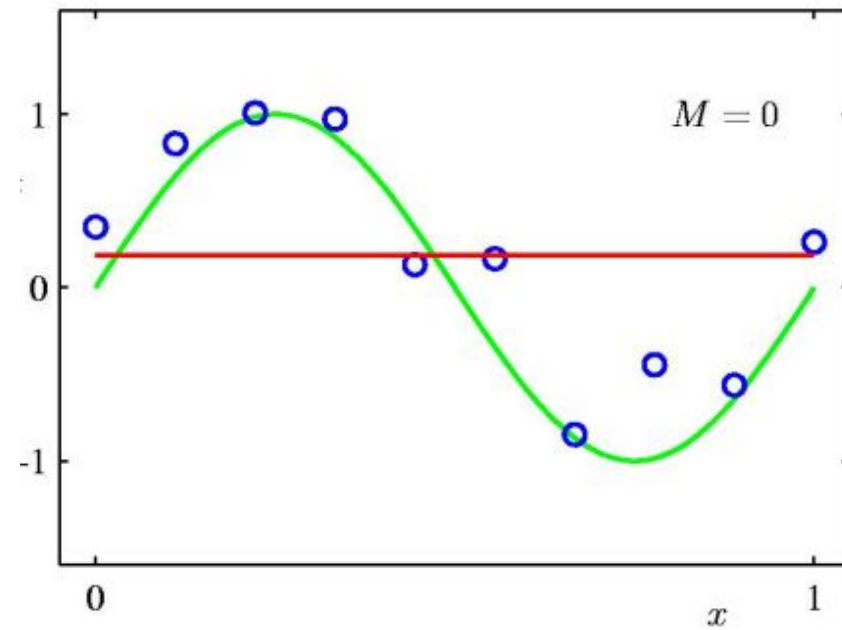
$$y(x, \mathbf{w}) = w_0 + \sum_{j=1}^M w_j x^j$$

where x^j is the j -th power of x

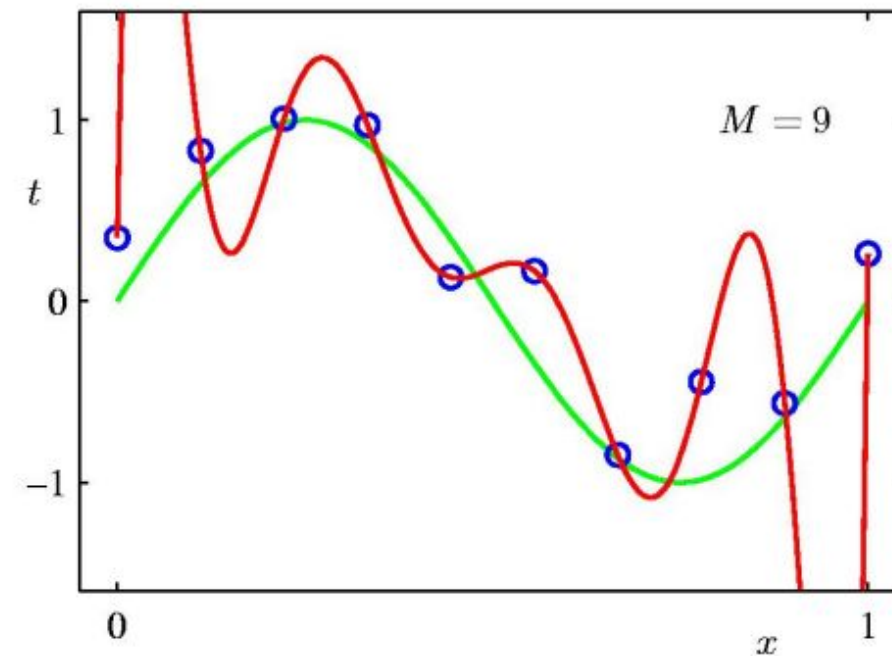
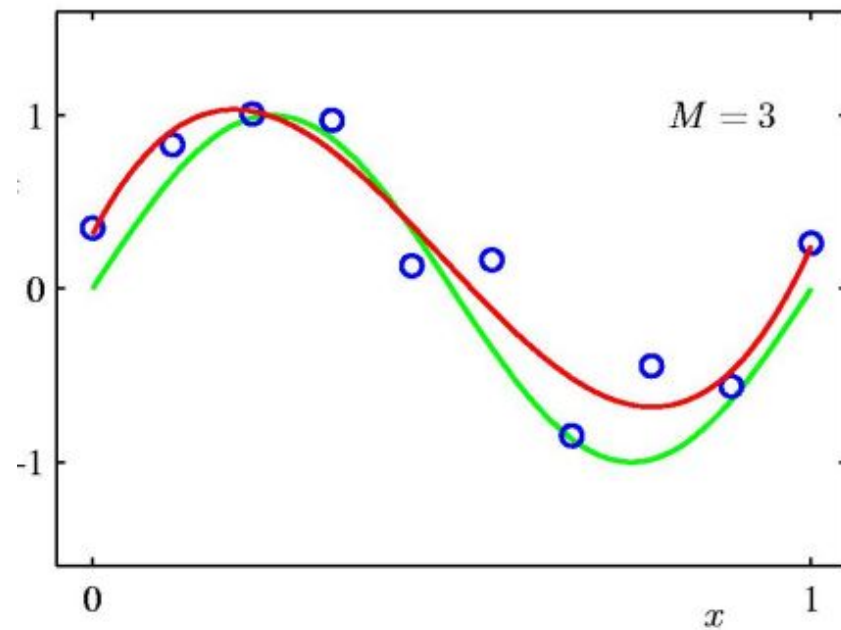
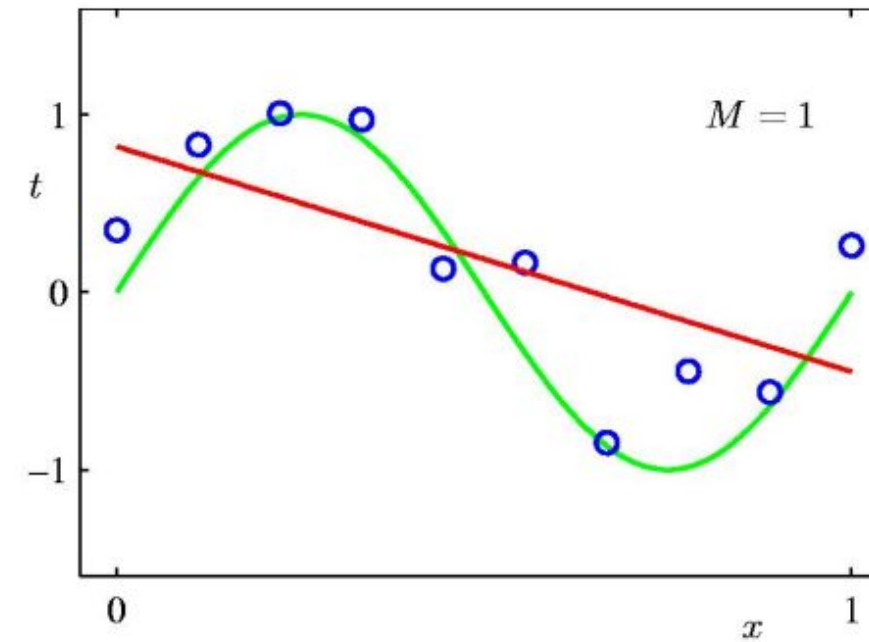
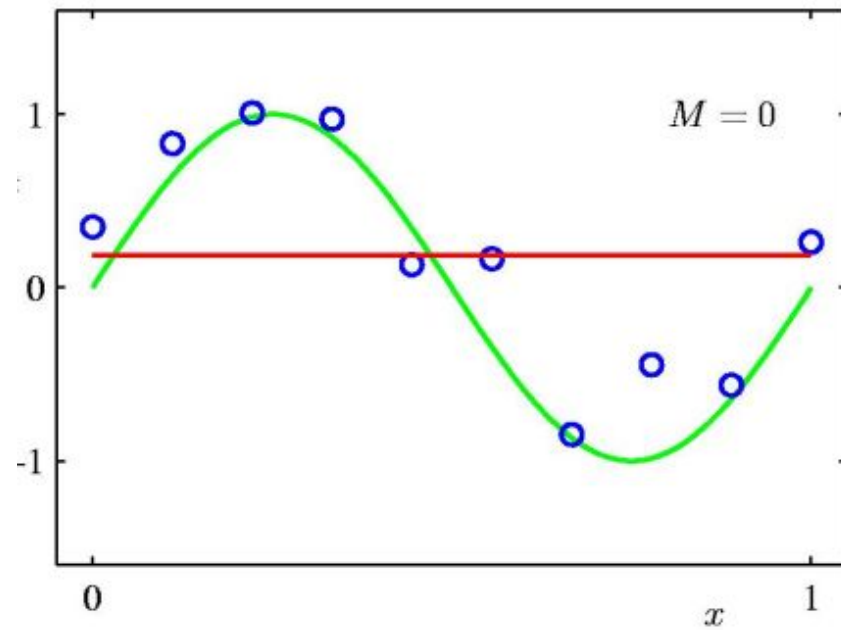
We can use the same approach to optimize for the **weights w**



Fitting a Polynomial



Fitting a Polynomial

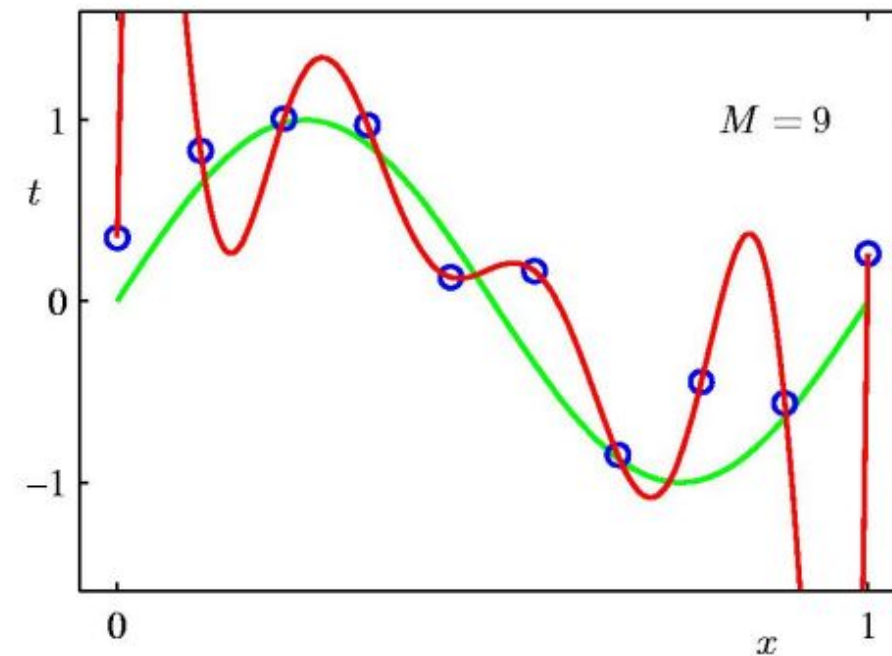
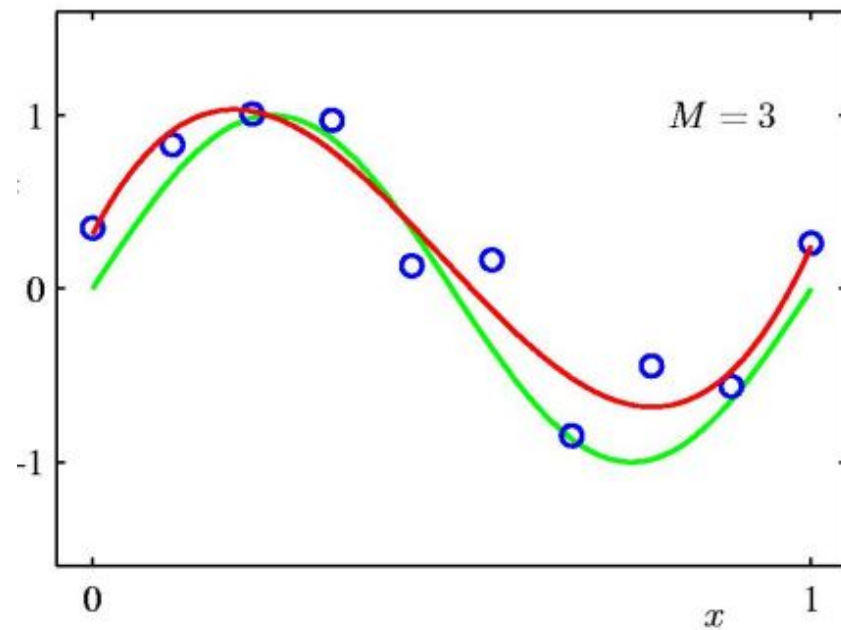
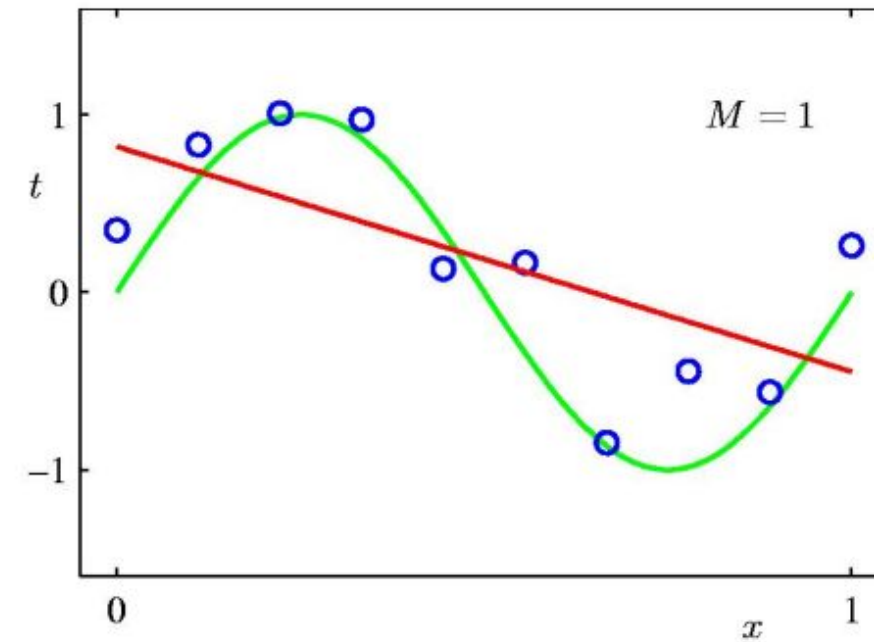
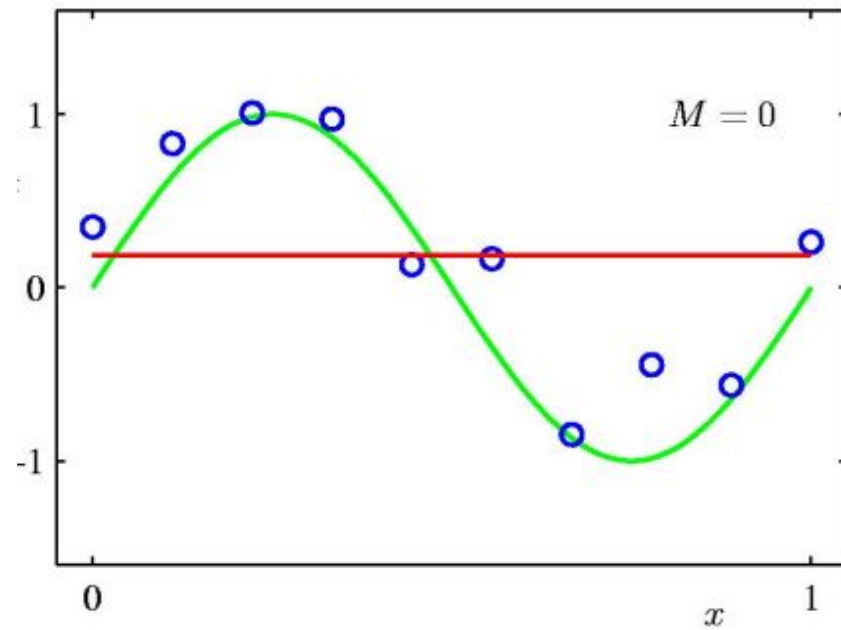


Error

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - f(x))^2$$

Fitting a Polynomial



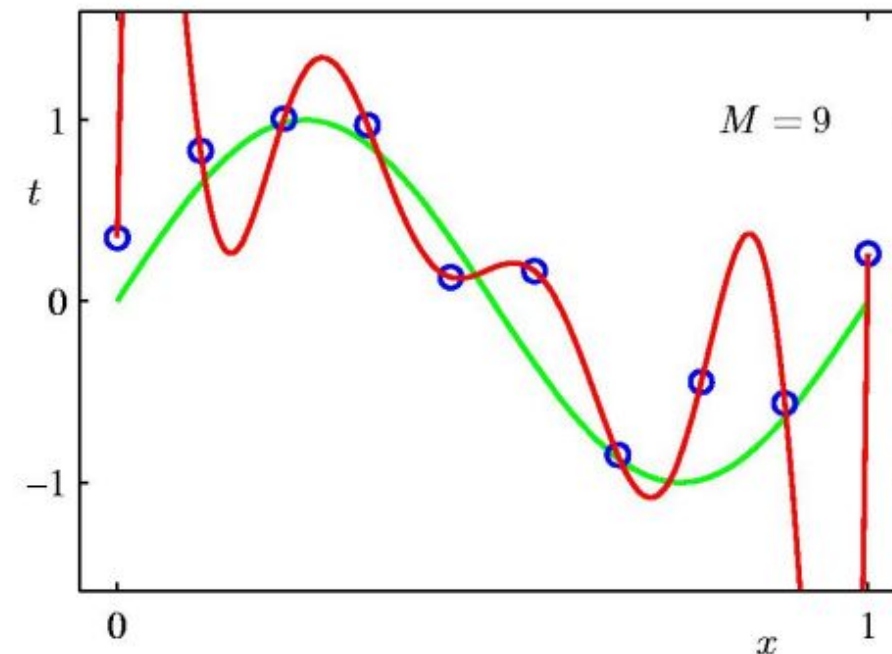
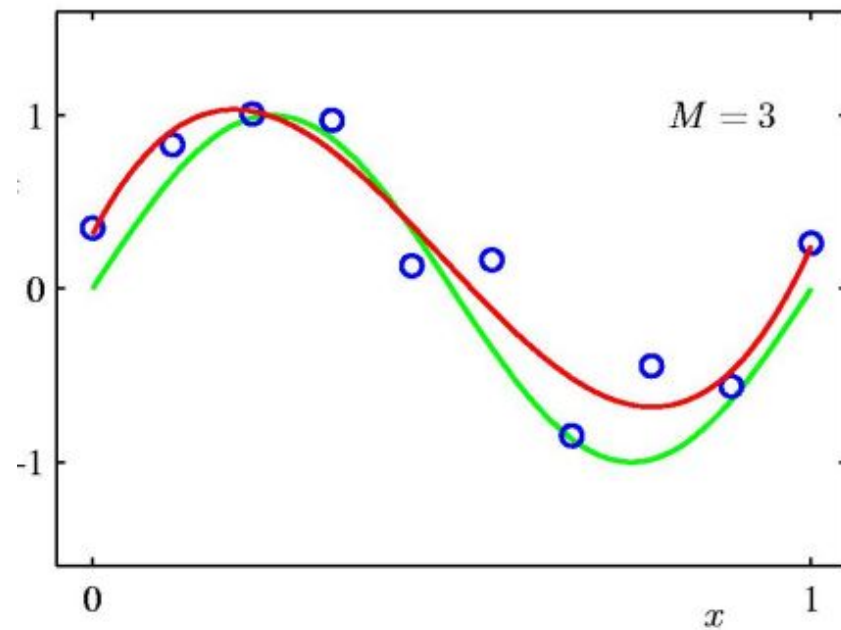
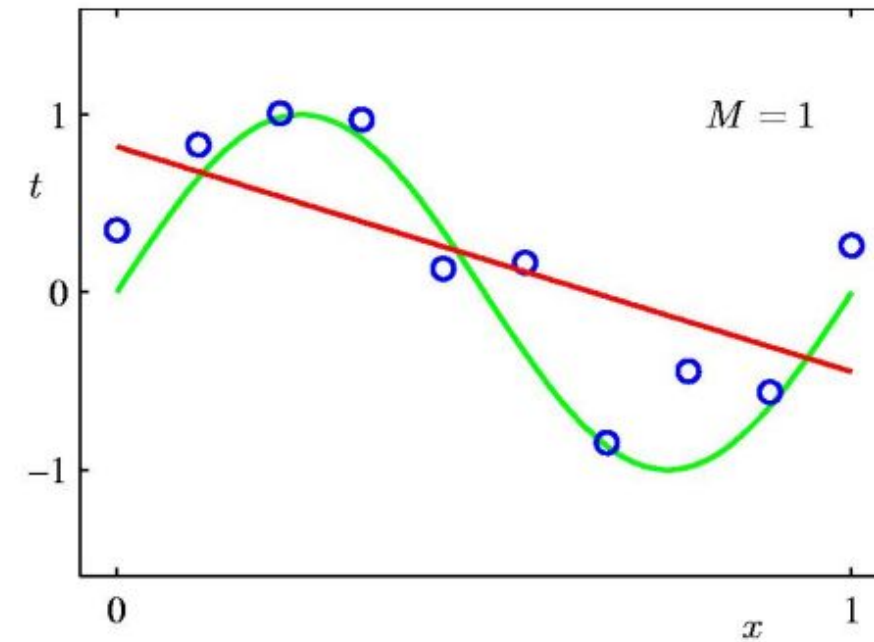
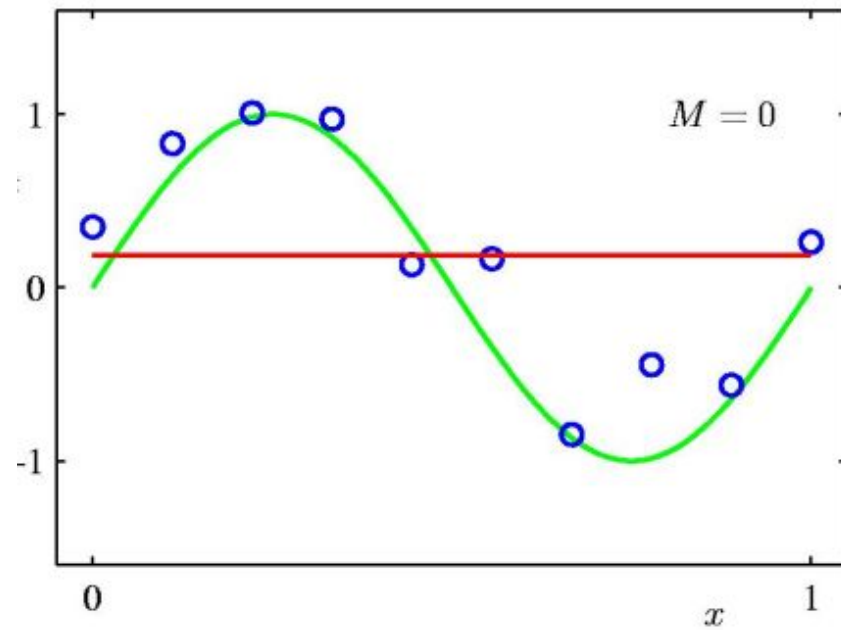
Error

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - f(x))^2$$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - (wx + b))^2$$

Fitting a Polynomial



Error

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - f(x))^2$$

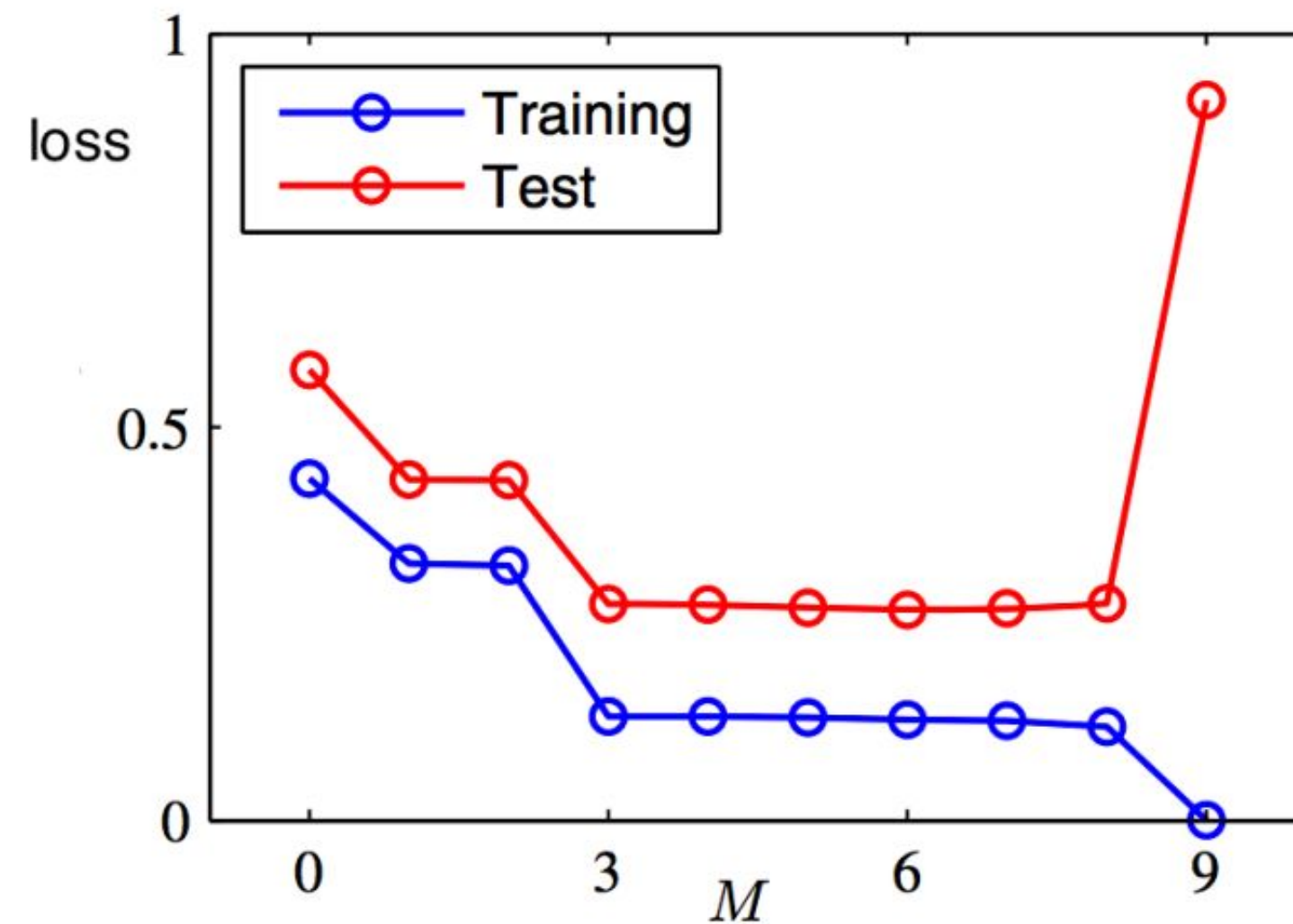
$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - (wx + b))^2$$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - (w_0x^2 + w_1x + b))^2$$

Generalization

Generalization = model's ability to predict the held out data

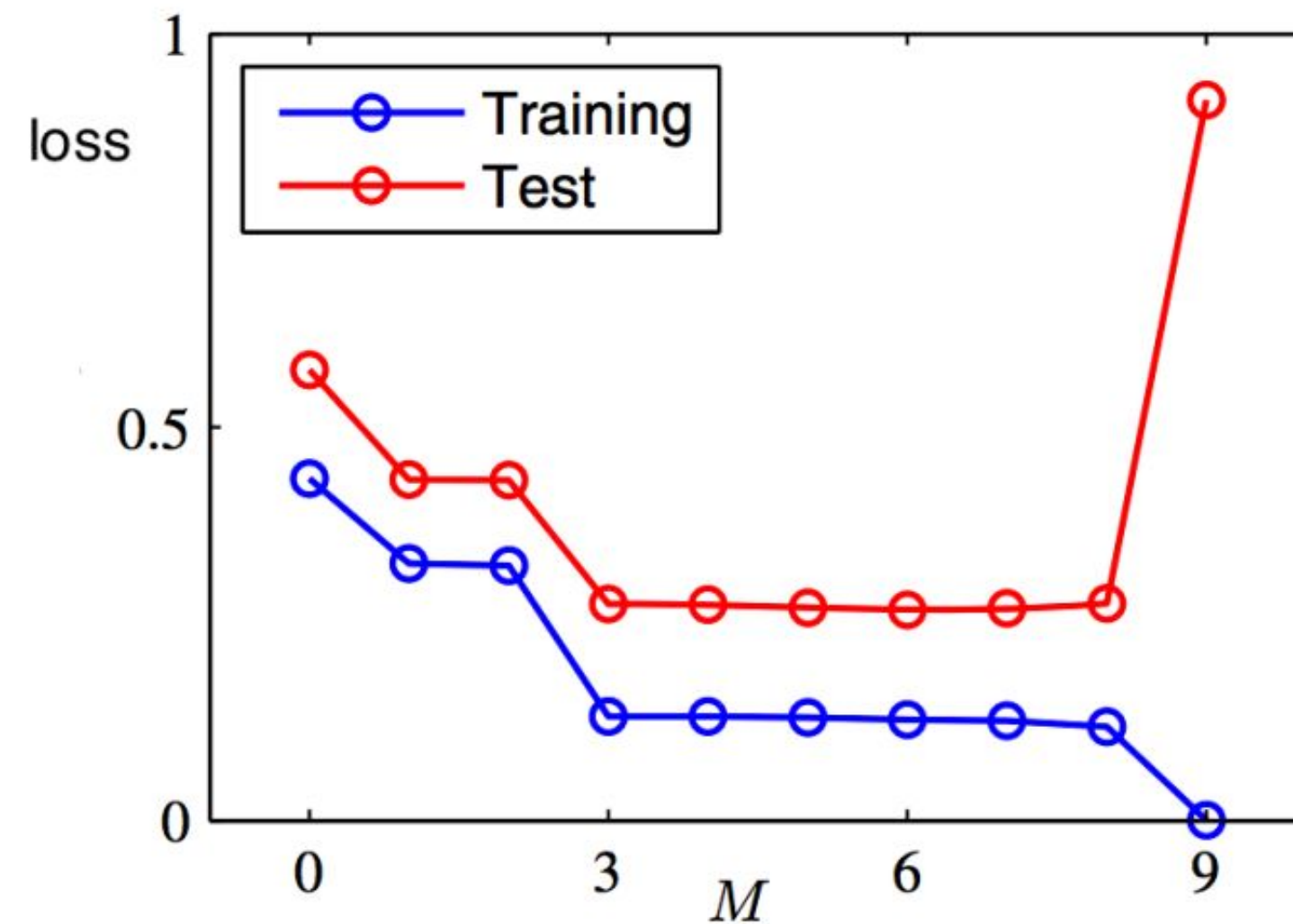
What is happening?



Generalization

Generalization = model's ability to predict the held out data

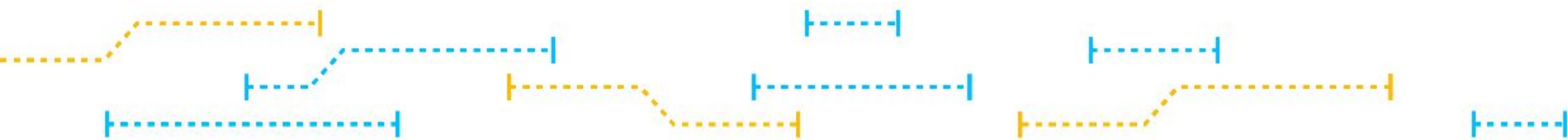
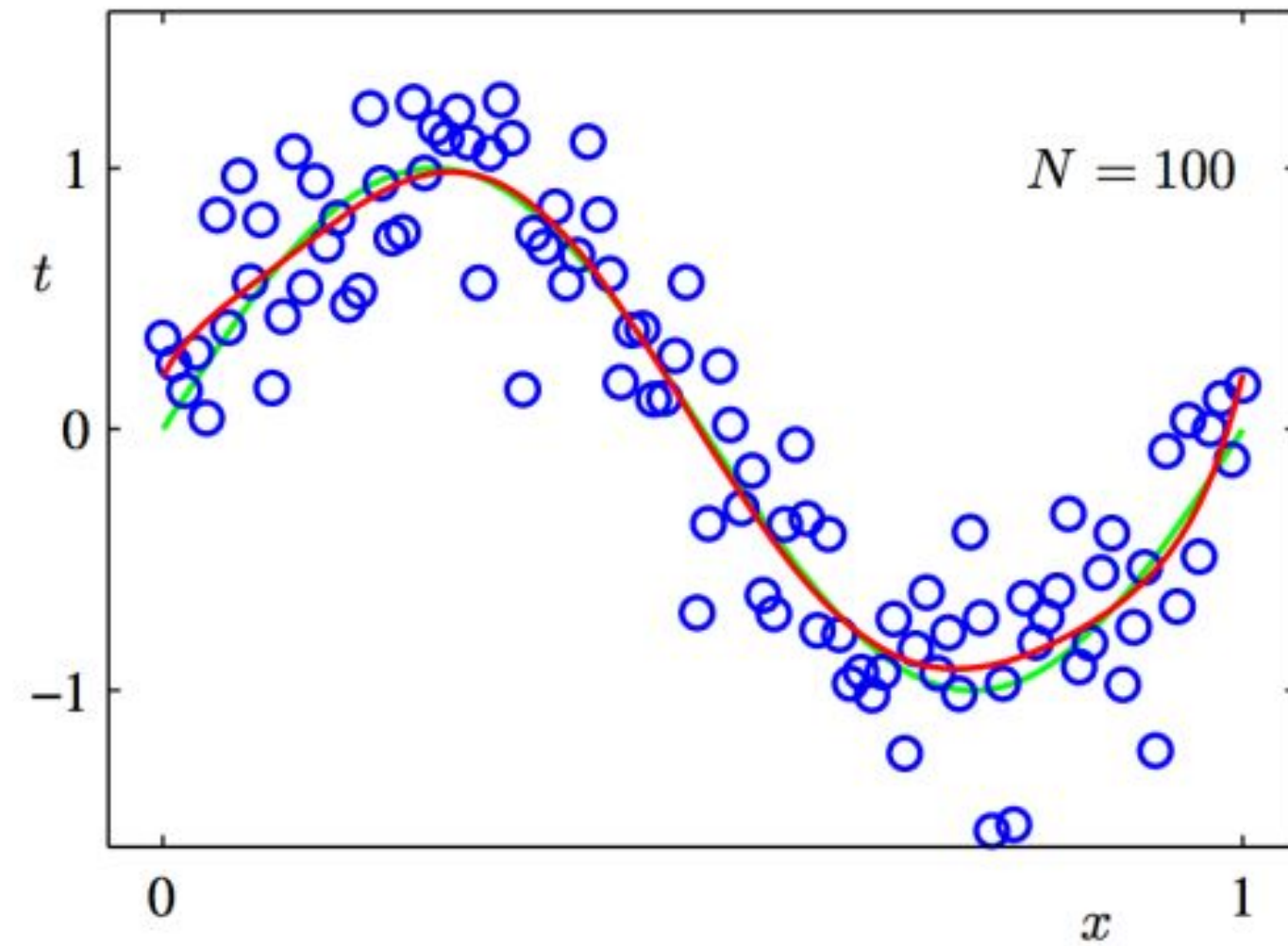
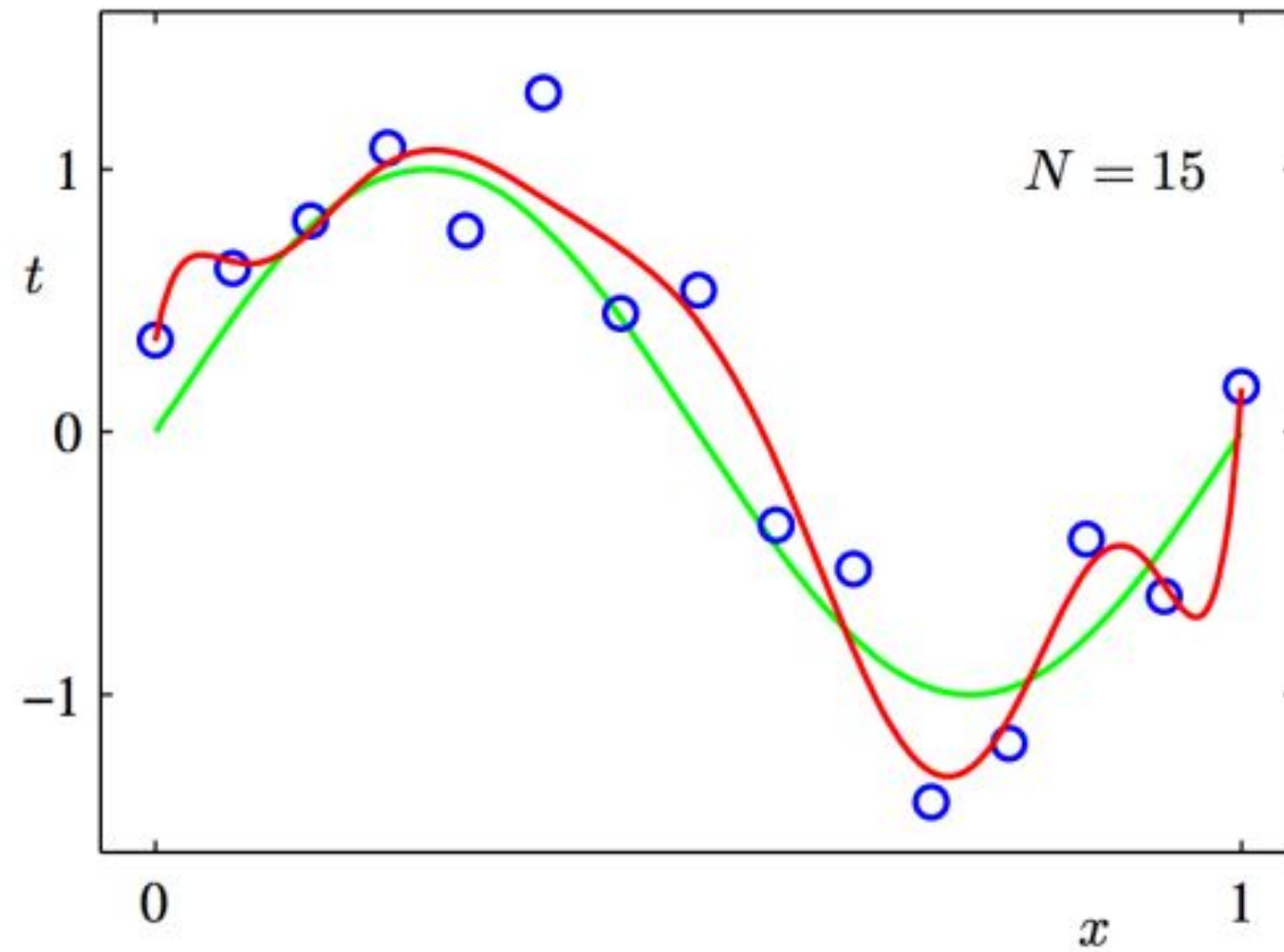
What is happening?



Our model with $M = 9$ **overfits** the data (it models also noise)

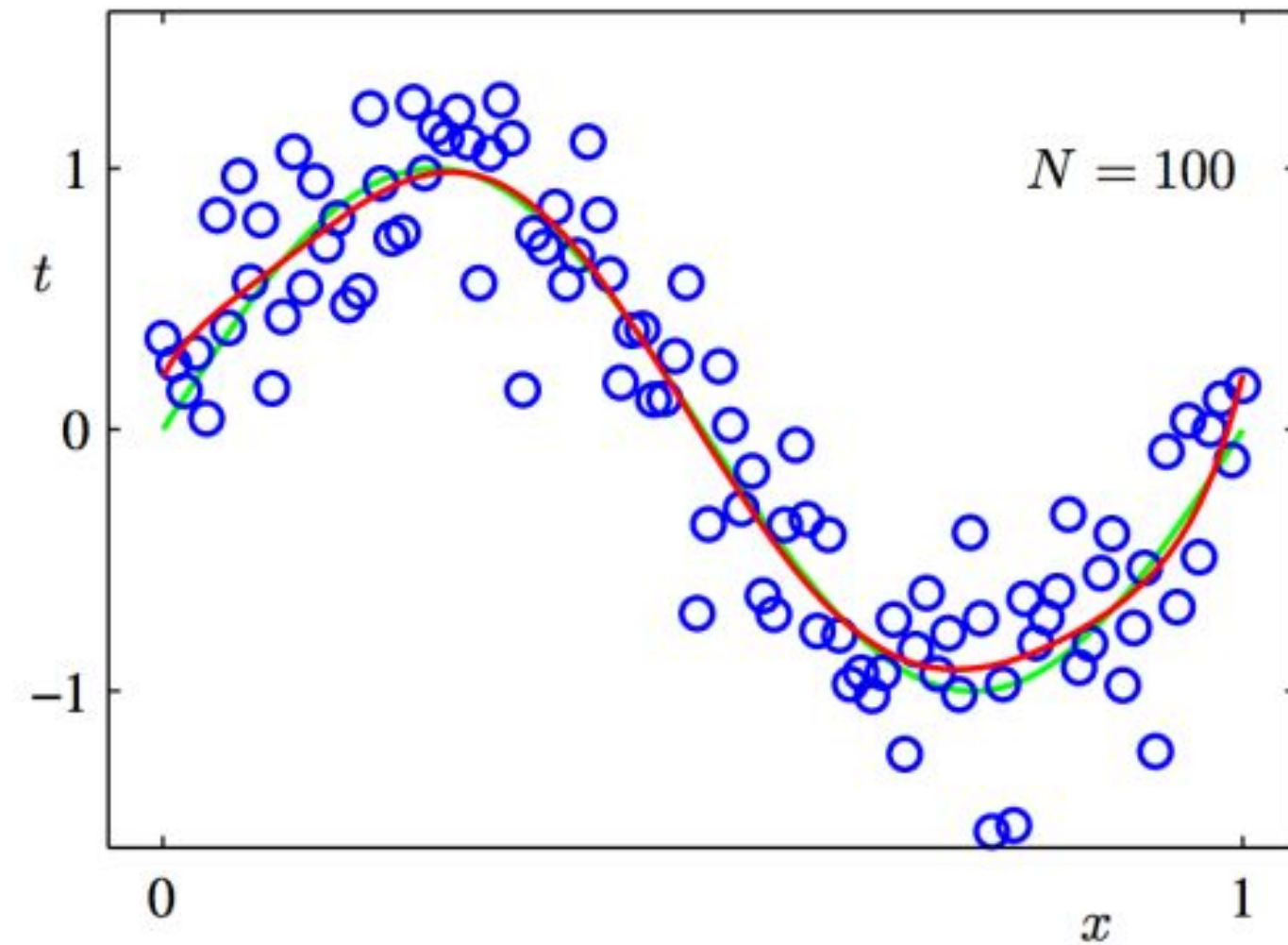
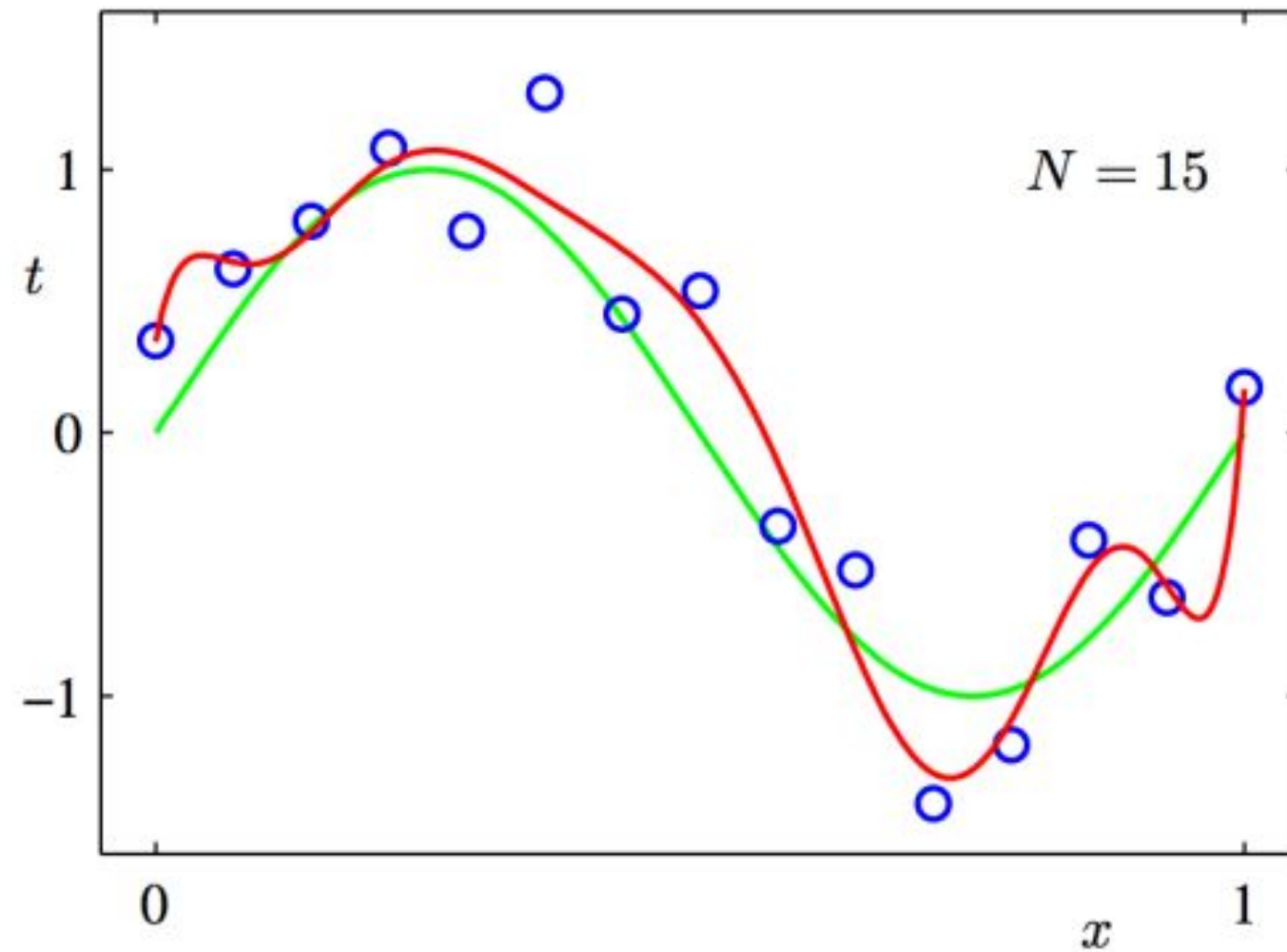
Generalization

What can we do?



Generalization

What can we do? Not a problem if we have lots of training examples



Generalization

Let's look at the **estimated weights** for **various M** in the case of fewer examples

	$M = 0$	$M = 1$	$M = 6$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43



Generalization

Let's look at the **estimated weights** for **various M** in the case of fewer examples

The weights are becoming huge to compensate for the noise

	$M = 0$	$M = 1$	$M = 6$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43





How can we deal with this?



3.



Regularization

L0 & L1 & L2

Regularization

Occam's Razor: prefer the simplest hypothesis

What does it mean for a hypothesis (or model) to be simple?

1. small number of features (**model selection**)
2. small number of "important" features (**shrinkage**)



Regularization

Given objective function: $J(\theta)$

Goal is to find

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} J(\theta) + \lambda r(\theta)$$

Key idea: Define regularizer $r(\theta)$

s.t. we tradeoff between fitting the data and keeping the model simple

Choose form of regularizer

Common choice: p-norm

$$r(\theta) = \|\theta\|_p = \left[\sum_{m=1}^M |\theta_m|^p \right]^{\left(\frac{1}{p}\right)}$$



Regularization

Given objective function: $J(\theta)$

Goal is to find

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} J(\theta) + \lambda r(\theta) \quad r(\theta) = \|\theta\|_p = \left[\sum_{m=1}^M |\theta_m|^p \right]^{\left(\frac{1}{p}\right)}$$

p	$r(\theta)$	yields parameters that are...	name	optimization notes
0	$\ \theta\ _0 = \sum \mathbb{1}(\theta_m \neq 0)$	zero values	L0 reg.	no good computational solutions
1	$\ \theta\ _1 = \sum \theta_m $	zero values	L1 reg.	subdifferentiable
2	$(\ \theta\ _2)^2 = \sum \theta_m^2$	small values	L2 reg.	differentiable

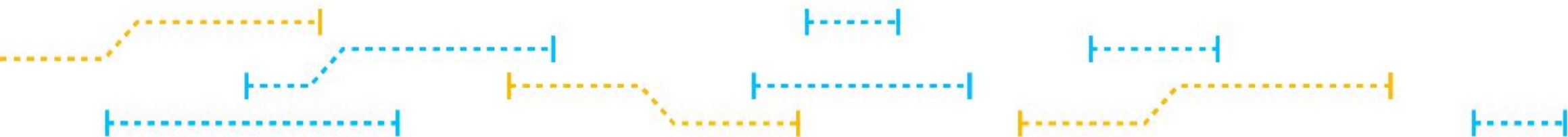


L2+L1 Regularization

L2 and L1 regularization can be **combined**

$$R(\mathbf{w}) = \lambda_2 \|\mathbf{w}\|^2 + \lambda_1 \|\mathbf{w}\|_1$$

- Also called **ElasticNet**
- Can work better than either type alone
- Can **adjust hyperparameters** to control which of the two penalties is more important



4.

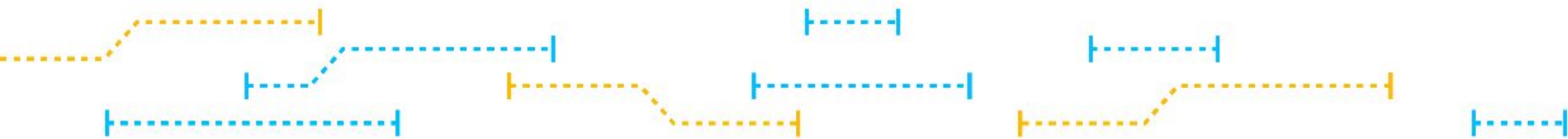
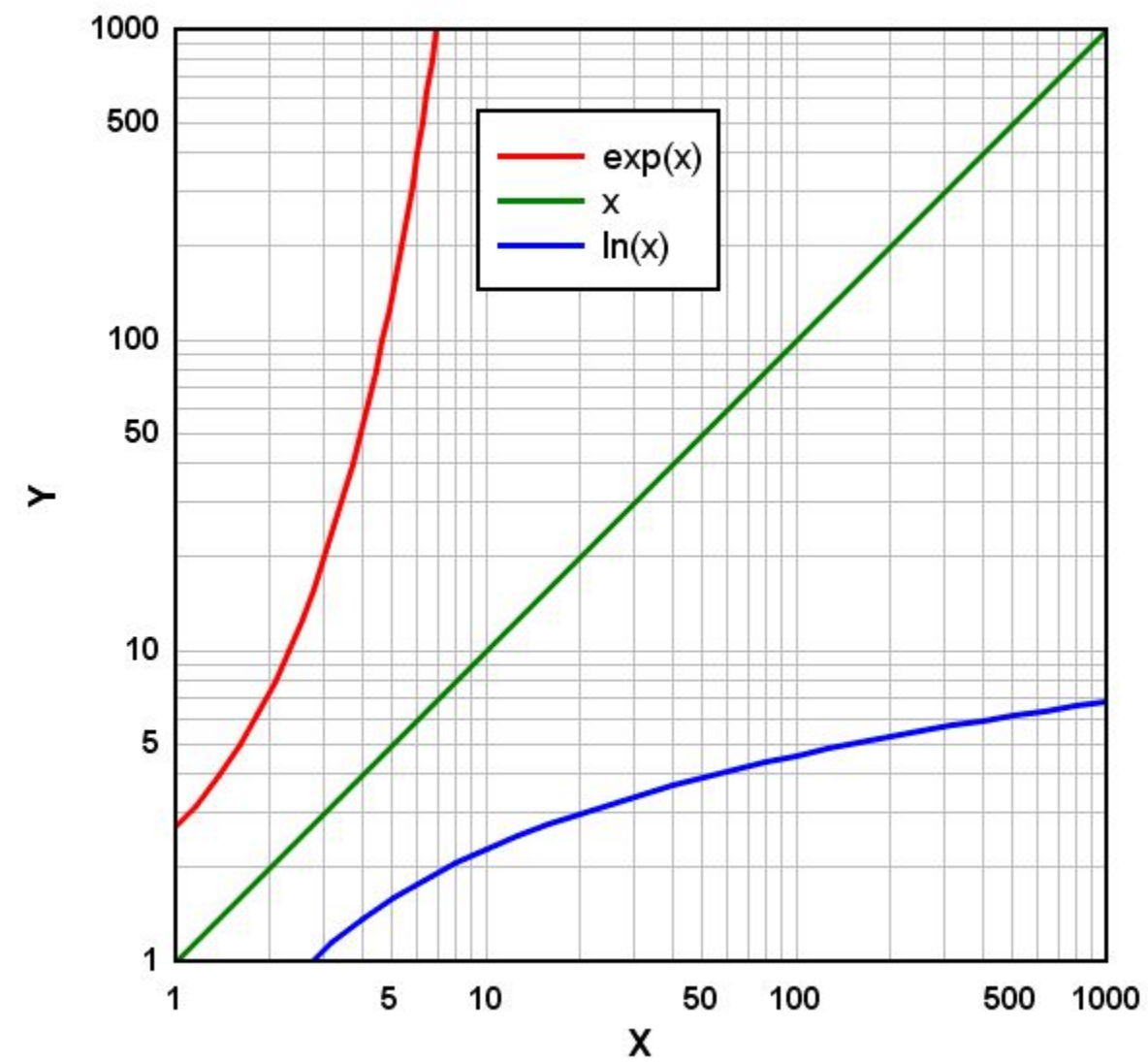


More non linear

Non polynomial

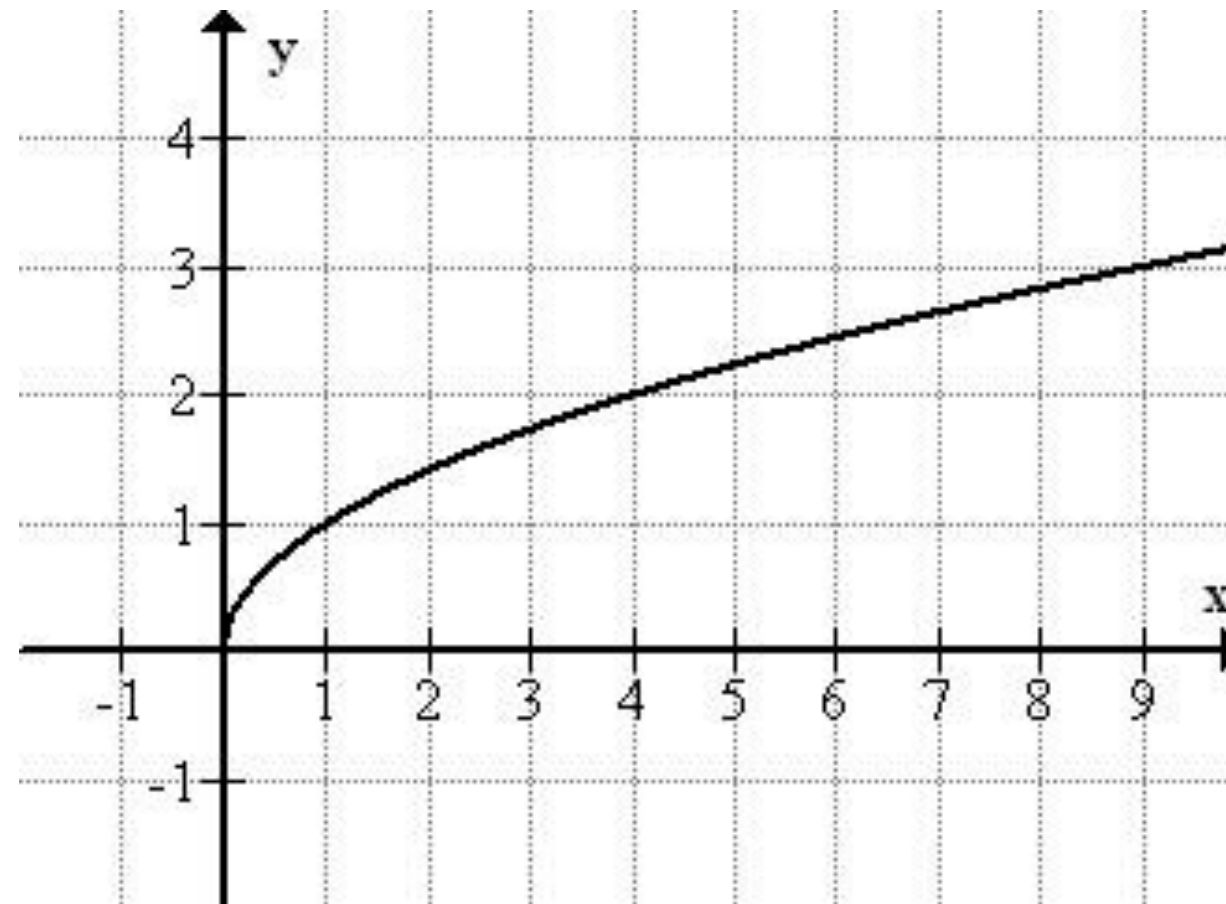
More non linear functions...

Logarithmic and exponential function



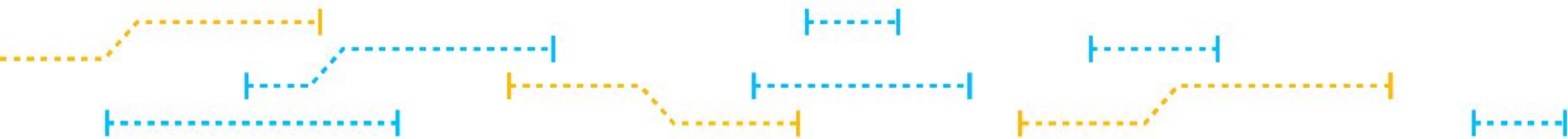
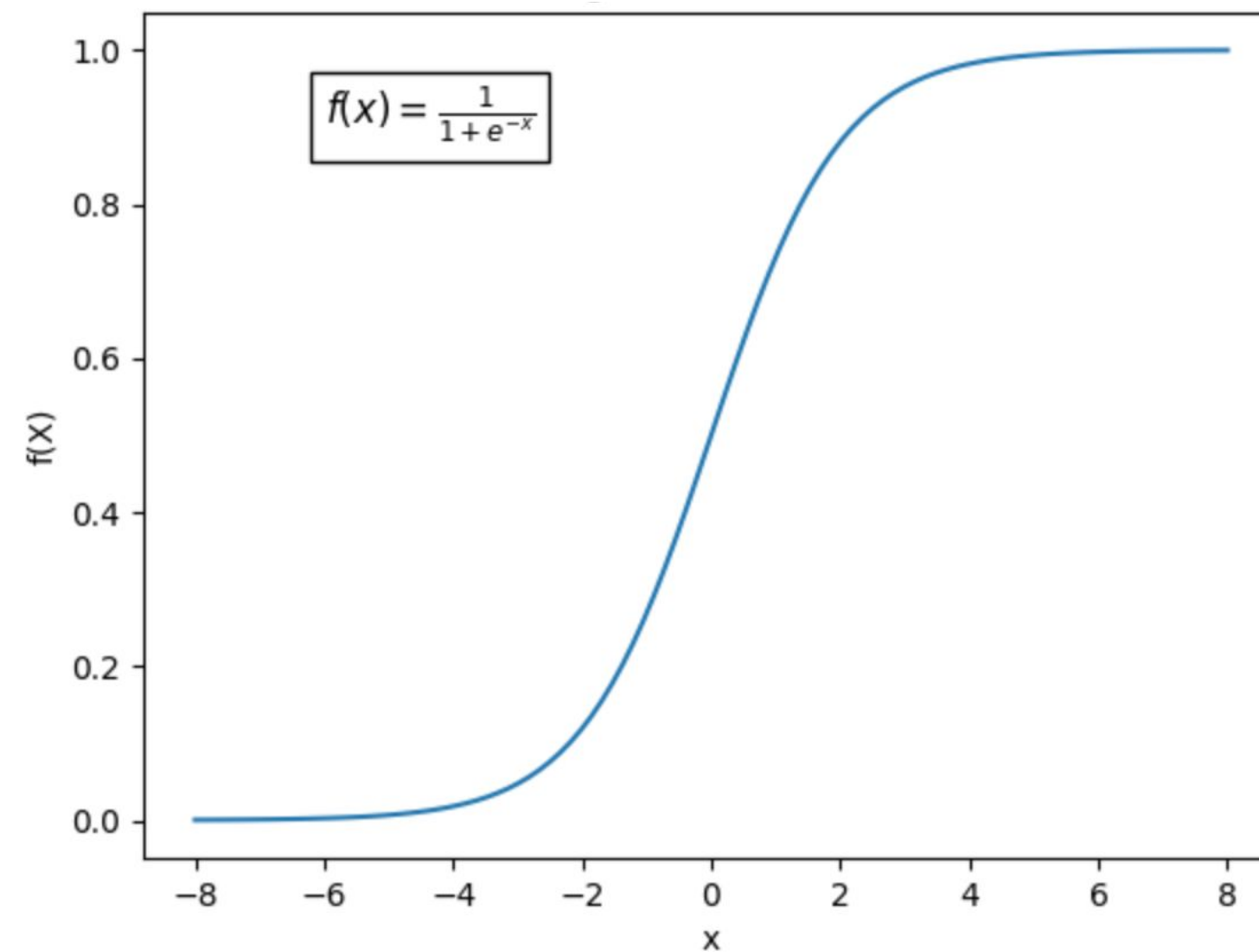
More non linear functions...

Square root function



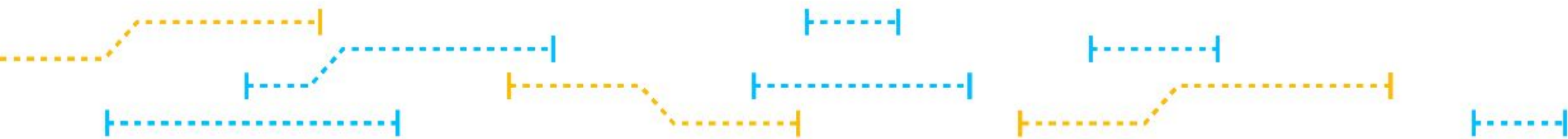
More non linear functions...

Sigmoid function



What about more dimensions?

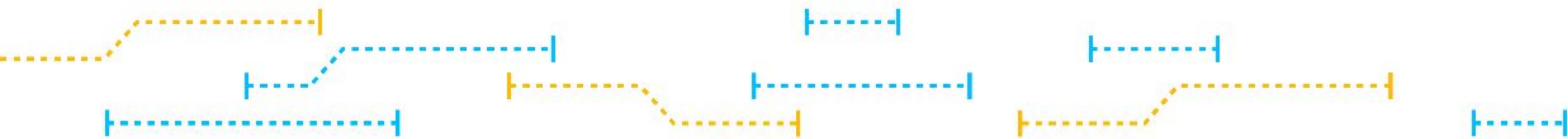
$$f(x) = w_0 + x_1 w_1 + x_2 w_2$$



What about more dimensions?

$$f(x) = w_0 + x_1 w_1 + x_2 w_2$$

$$f(x) = w_0 + x_1 w_1 + x_1^2 w_2 + x_2 w_3$$

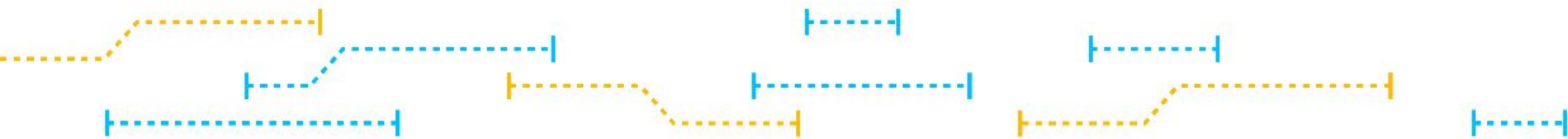


What about more dimensions?

$$f(x) = w_0 + x_1 w_1 + x_2 w_2$$

$$f(x) = w_0 + x_1 w_1 + x_1^2 w_2 + x_2 w_3$$

$$f(x) = w_0 + x_1 w_1 + x_1^2 w_2 + x_1^3 w_3 + x_2^2 w_4 + x_2 w_5$$



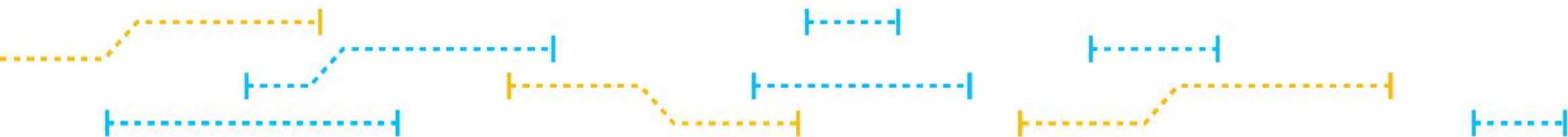
What about more dimensions?

$$f(x) = w_0 + x_1 w_1 + x_2 w_2$$

$$f(x) = w_0 + x_1 w_1 + x_1^2 w_2 + x_2 w_3$$

$$f(x) = w_0 + x_1 w_1 + x_1^2 w_2 + x_1^3 w_3 + x_2^2 w_4 + x_2 w_5$$

$$f(x) = w_0 + x_1 w_1 + x_1^2 w_2 + x_1^3 w_3 + x_2^2 w_4 + x_2 w_5 + x_1 x_2 w_6$$



What about more dimensions?

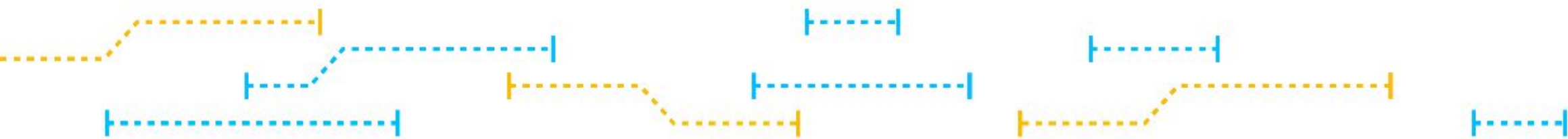
$$f(x) = w_0 + x_1 w_1 + x_2 w_2$$

$$f(x) = w_0 + x_1 w_1 + x_1^2 w_2 + x_2 w_3$$

$$f(x) = w_0 + x_1 w_1 + x_1^2 w_2 + x_1^3 w_3 + x_2^2 w_4 + x_2 w_5$$

$$f(x) = w_0 + x_1 w_1 + x_1^2 w_2 + x_1^3 w_3 + x_2^2 w_4 + x_2 w_5 + x_1 x_2 w_6$$

$$f(x) = w_0 + \exp(x_1) w_1 + x_2 w_2$$



What about more dimensions?

$$f(x) = w_0 + x_1 w_1 + x_2 w_2$$

$$f(x) = w_0 + x_1 w_1 + x_1^2 w_2 + x_2 w_3$$

$$f(x) = w_0 + x_1 w_1 + x_1^2 w_2 + x_1^3 w_3 + x_2^2 w_4 + x_2 w_5$$

$$f(x) = w_0 + x_1 w_1 + x_1^2 w_2 + x_1^3 w_3 + x_2^2 w_4 + x_2 w_5 + x_1 x_2 w_6$$

$$f(x) = w_0 + \exp(x_1) w_1 + x_2 w_2$$

$$f(x) = w_0 + \exp(x_1) w_1 + x_2^2 w_2 + x_2 w_3$$



What about more dimensions?

$$f(x) = w_0 + x_1 w_1 + x_2 w_2$$

$$f(x) = w_0 + x_1 w_1 + x_1^2 w_2 + x_2 w_3$$

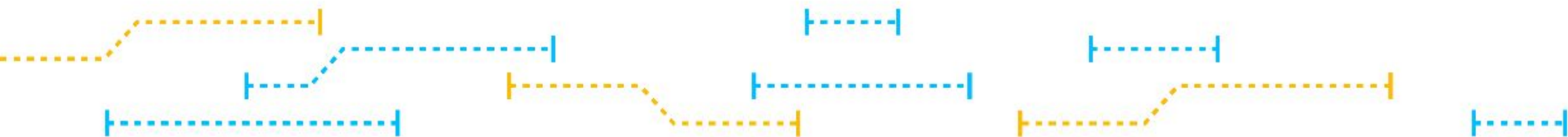
$$f(x) = w_0 + x_1 w_1 + x_1^2 w_2 + x_1^3 w_3 + x_2^2 w_4 + x_2 w_5$$

$$f(x) = w_0 + x_1 w_1 + x_1^2 w_2 + x_1^3 w_3 + x_2^2 w_4 + x_2 w_5 + x_1 x_2 w_6$$

$$f(x) = w_0 + \exp(x_1) w_1 + x_2 w_2$$

$$f(x) = w_0 + \exp(x_1) w_1 + x_2^2 w_2 + x_2 w_3$$

•
•
•



5.



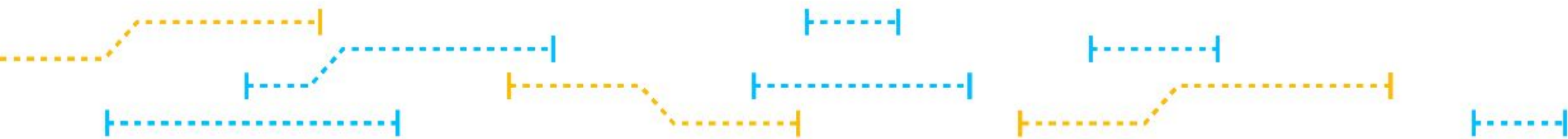
Non Linear Regression

Example

https://colab.research.google.com/drive/1-ZJqtjIM_PwAdN2vMmYzDRDv-DH9ILdj?usp=sharing

More Exercises?

MACHINE LEARNING WORKOUT: Chapter 5 and 6



Kahoot Time!

