



#### 10-301/601 Introduction to Machine Learning

Machine Learning Department School of Computer Science Carnegie Mellon University

## (Linear Models) + Feature Engineering + Regularization

Matt Gormley Lecture 10 Sep. 30, 2022

#### Q&A

- I think you graded seven different questions incorrectly. Should I put them all in on Gradescope regrade request and submit that?
- A: Please, no. I'd encourage you to watch this tutorial video from Gradescope so you know how to use this important tool.

https://help.gradescope.com/article/8hchz9h8wh-student-regrade-request

#### Reminders

- Practice Problems: Exam 1
- Exam 1
  - Tue, Oct 4, 6:30pm 8:30pm
  - see Piazza for details
- Homework 4: Logistic Regression
  - Out: Tue, Oct 4
  - Due: Thu, Oct 13 at 11:59pm

Linear Models

### PERCEPTRON, LINEAR REGRESSION, AND LOGISTIC REGRESSION

#### Why is it not "Logistic Classification"?

#### Whiteboard

- Conceptual Change: 2D classification in 3D
- Why is it called Logistic Regression and not Logistic Classification?

#### Q1

#### Matching Game

6.

#### **Question:**

Match the Algorithm to its Update Rule

1. SGD for Logistic Regression

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = p(\hat{y}|\hat{x}) = \sigma(\hat{\boldsymbol{\theta}}^{\mathsf{T}}\hat{\mathbf{x}})$$

2. Least Mean Squares

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}$$

3. Perceptron

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = \operatorname{sign}(\boldsymbol{\theta}^T \mathbf{x})$$

4.  $\theta_k \leftarrow \theta_k + (h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)})$ 

$$\theta_k \leftarrow \theta_k + \frac{1}{1 + \exp \lambda (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)})}$$

 $\theta_k \leftarrow \theta_k + \lambda (h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)}) x_k^{(i)}$ 

#### **Answer:**

$$C. 1=6, 2=4, 3=4$$



#### SGD for Logistic Regression

#### **Question:**

Which of the following is a correct description of SGD for Logistic Regression?

#### **Answer:**

At each step (i.e. iteration) of SGD for Logistic Regression we...

- A. (1) compute the gradient of the log-likelihood for all examples (2) update all the parameters using the gradient
- B. (1) ask Matt for a description of SGD for Logistic Regression, (2) write it down, (3) report that answer loxic
- C. (1) compute the gradient of the log-likelihood for all examples (2) randomly pick an example (3) update only the parameters for that example
- D. (1) randomly pick a parameter, (2) compute the partial derivative of the log-likelihood with respect to that parameter, (3) update that parameter for all examples
- (1) randomly pick an example, (2) compute the gradient of the log-likelihood for that example, (3) update all the parameters using that gradient
  - F. (1) randomly pick a parameter and an example, (2) compute the gradient of the log-likelihood for that example with respect to that parameter, (3) update that parameter using that gradient



#### **Gradient Descent**

#### Algorithm 1 Gradient Descent

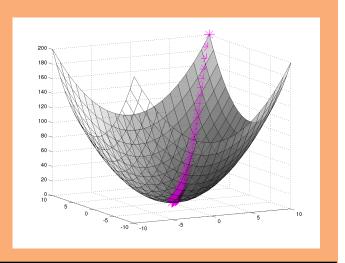
1: **procedure**  $GD(\mathcal{D}, \boldsymbol{\theta}^{(0)})$ 

2: 
$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}^{(0)}$$

3: while not converged do

4: 
$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \gamma \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

5: return  $\theta$ 



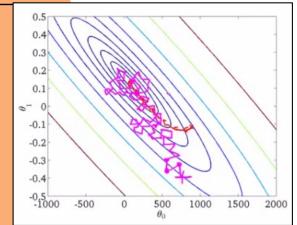
In order to apply GD to Logistic Regression all we need is the **gradient** of the objective function (i.e. vector of partial derivatives).

$$abla_{m{ heta}} J(m{ heta}) = egin{bmatrix} rac{d heta_1}{d heta_2} J(m{ heta}) \ rac{d}{d heta_2} J(m{ heta}) \ rac{d}{d heta_M} J(m{ heta}) \end{bmatrix}$$

## Stochastic Gradient Descent (SGD)

#### Algorithm 1 Stochastic Gradient Descent (SGD)

```
1: \operatorname{procedure} \operatorname{SGD}(\mathcal{D}, \boldsymbol{\theta}^{(0)})
2: \boldsymbol{\theta} \leftarrow \boldsymbol{\theta}^{(0)}
3: \operatorname{while} not converged \operatorname{do}
4: \operatorname{for} i \in \operatorname{shuffle}(\{1, 2, \dots, N\}) \operatorname{do}
5: \boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \gamma \nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta})
```



We can also apply SGD to solve the MCLE problem for Logistic Regression.

We need a per-example objective:

return  $\theta$ 

6:

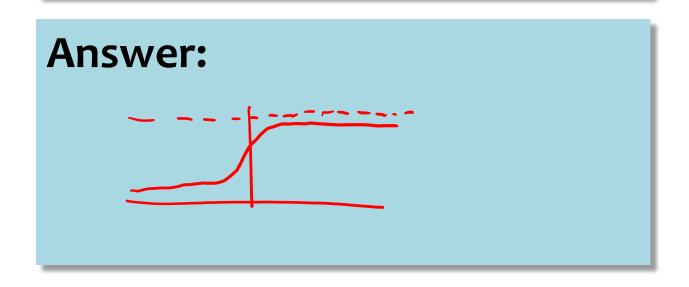
Let 
$$J(\boldsymbol{\theta}) = \sum_{i=1}^{N} J^{(i)}(\boldsymbol{\theta})$$
 where  $J^{(i)}(\boldsymbol{\theta}) = -\log p_{\boldsymbol{\theta}}(y^i|\mathbf{x}^i)$ .

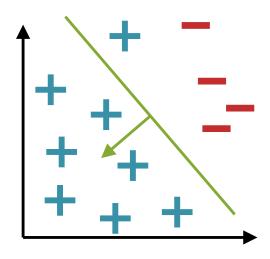
Logistic Regression vs. Perceptron uestion:

A= toxic B= Toxic C= False 75%

**Question:** 

True or False: Just like Perceptron, one step (i.e. iteration) of SGD for Logistic Regression will result in a change to the parameters only if the current example is incorrectly classified.





#### **BAYES OPTIMAL CLASSIFIER**

#### Bayes Optimal Classifier

Suppose you knew the distribution p\*(y | x) or had a good approximation to it.

#### **Question:**

How would you design a function y = h(x) to predict a single label?  $\in \{0, 13\}$ 

#### **Answer:**

You'd use the Bayes optimal classifier!

approximates c (x)

#### **Probabilistic Learning**

Today, we assume that our output is **sampled** from a conditional **probability distribution**:

$$\mathbf{x}^{(i)} \sim p^*(\cdot)$$

$$y^{(i)} \sim p^*(\cdot|\mathbf{x}^{(i)})$$

Our goal is to learn a probability distribution p(y|x) that best approximates  $p^*(y|x)$ 

#### Bayes Optimal Classifier

#### Whiteboard

- Bayes Optimal Classifier
- Reducible / irreducible error
- Ex: Bayes Optimal Classifier for o/1 Loss

### OPTIMIZATION METHOD #4: MINI-BATCH SGD

#### Mini-Batch SGD

#### Gradient Descent:

Compute true gradient exactly from all N examples

#### Stochastic Gradient Descent (SGD):

Approximate true gradient by the gradient of one randomly chosen example

#### Mini-Batch SGD:

Approximate true gradient by the average gradient of K randomly chosen examples

#### Mini-Batch SGD

#### while not converged: $\theta \leftarrow \theta - \gamma \mathbf{g}$

#### Three variants of first-order optimization:

Gradient Descent: 
$$\mathbf{g} = \nabla J(\pmb{\theta}) = \frac{1}{N} \sum_{i=1}^N \nabla J^{(i)}(\pmb{\theta})$$

SGD: 
$$\mathbf{g} = \nabla J^{(i)}(\boldsymbol{\theta})$$
 where  $i$  sampled uniformly

Mini-batch SGD: 
$$\mathbf{g} = \frac{1}{S} \sum_{s=1}^S \nabla J^{(i_s)}(\pmb{\theta})$$
 where  $i_s$  sampled uniformly  $\forall s$ 

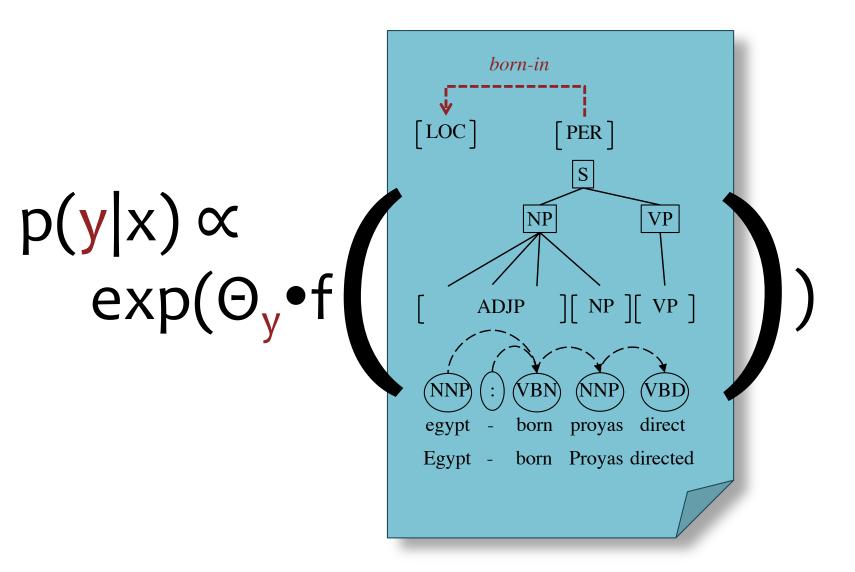
#### Logistic Regression Objectives

#### You should be able to...

- Apply the principle of maximum likelihood estimation (MLE) to learn the parameters of a probabilistic model
- Given a discriminative probabilistic model, derive the conditional log-likelihood, its gradient, and the corresponding Bayes Classifier
- Explain the practical reasons why we work with the log of the likelihood
- Implement logistic regression for binary or multiclass classification
- Prove that the decision boundary of binary logistic regression is linear
- For linear regression, show that the parameters which minimize squared error are equivalent to those that maximize conditional likelihood

#### FEATURE ENGINEERING

#### Handcrafted Features

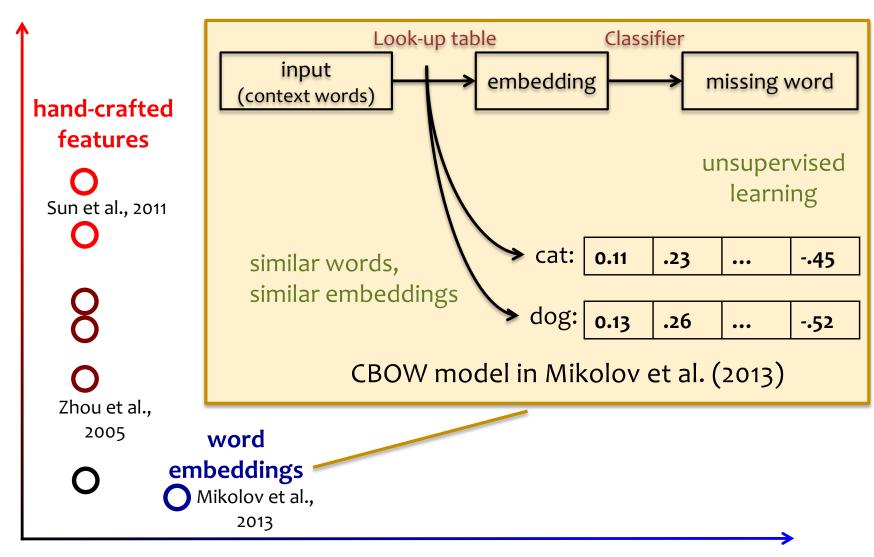


# Feature Engineering

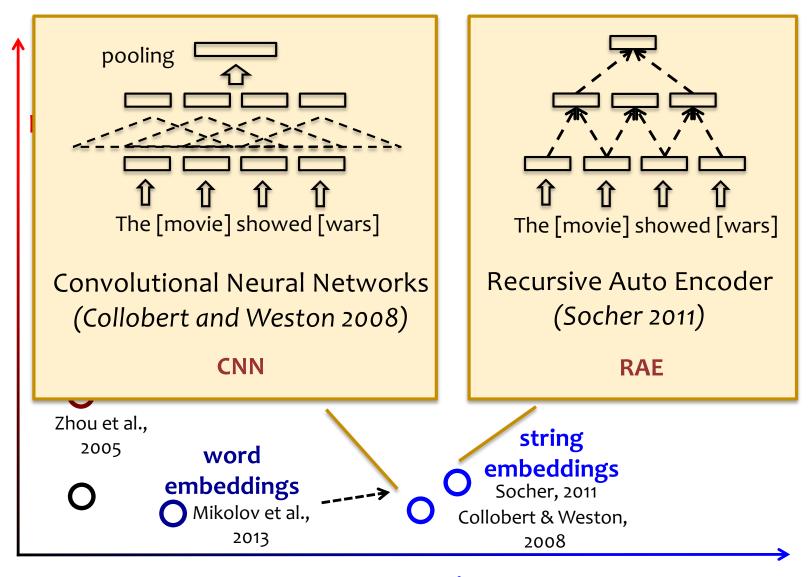
#### Where do features come from?

First word before M1 Second word before M1 hand-crafted Bag-of-words in M1 features Head word of M1 Other word in between First word after M2 Sun et al., 2011 Second word after M2 Bag-of-words in M2 Head word of M2 Bigrams in between Words on dependency path Country name list Personal relative triggers Personal title list Zhou et al., WordNet Tags 2005 Heads of chunks in between Path of phrase labels Combination of entity types

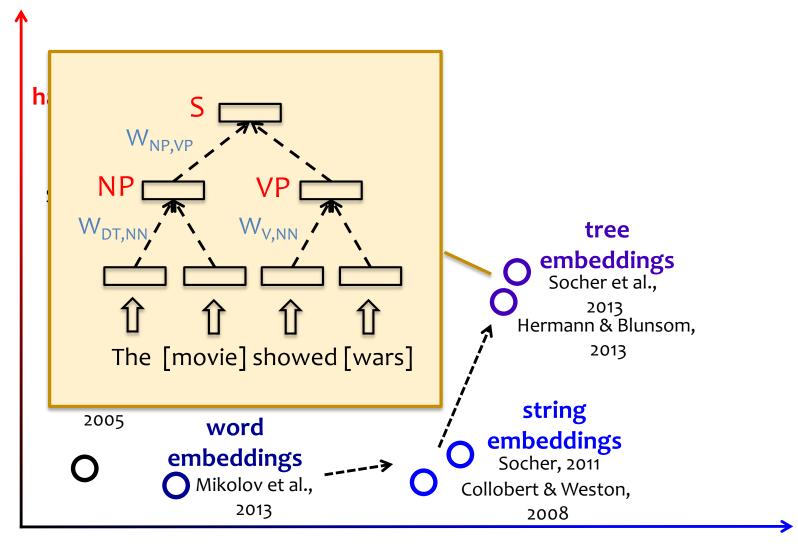
# Feature Engineering



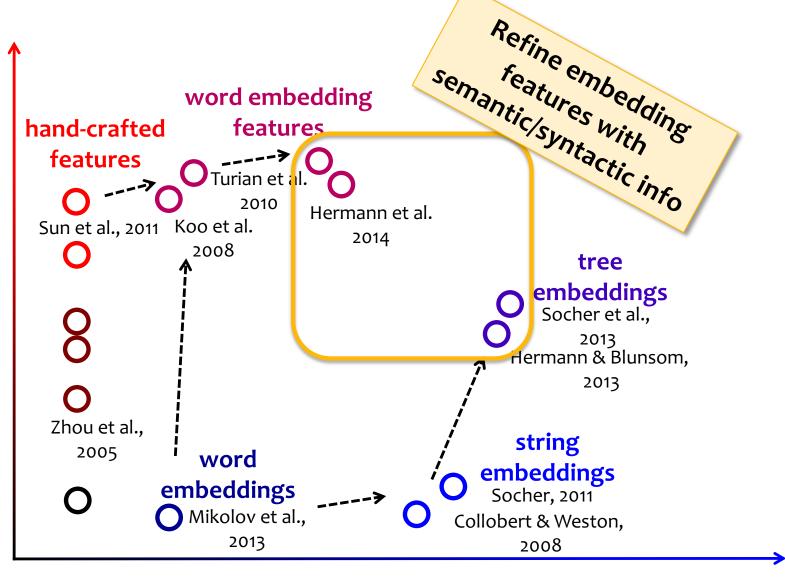
Feature Learning



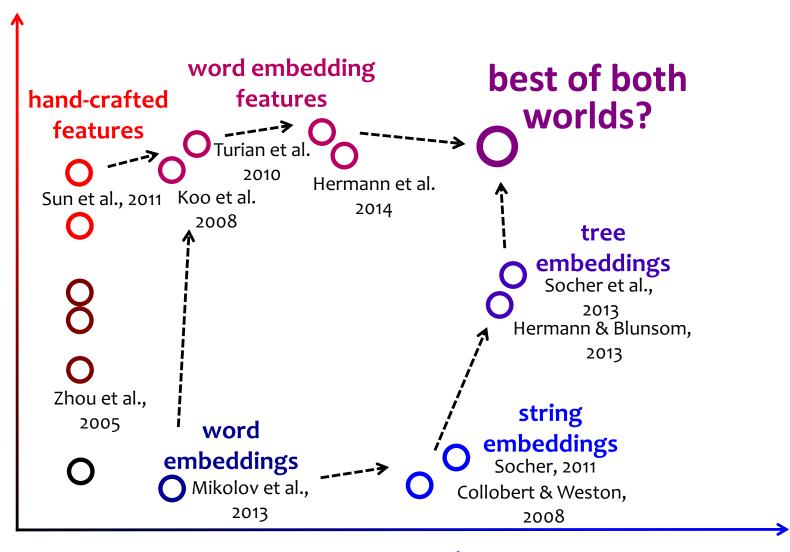
Feature Learning



Feature Learning



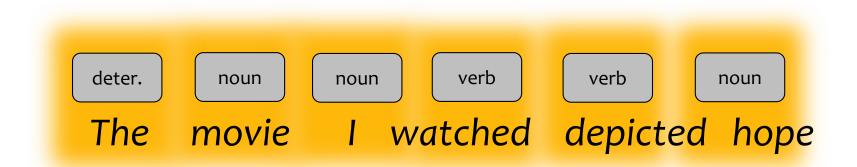
Feature Learning



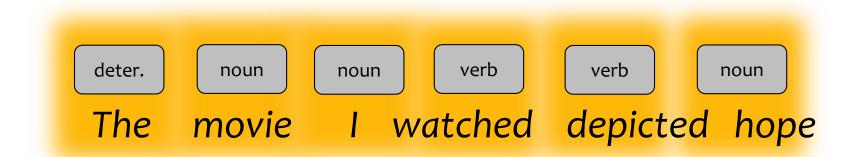
Feature Learning

Suppose you build a logistic regression model to predict a part-of-speech (POS) tag for each word in a sentence.

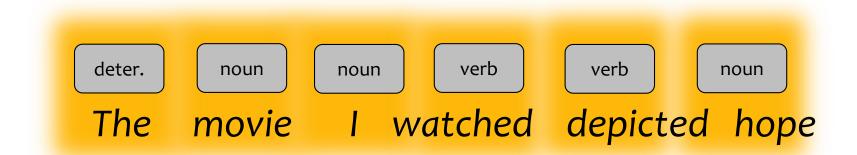
#### What features should you use?



#### **Per-word Features:**



#### **Context Features:**



#### **Context Features:**

 $w_{i} == "I"$   $w_{i+1} == "I"$   $w_{i-1} == "I"$   $w_{i+2} == "I"$   $w_{i-2} == "I"$ 

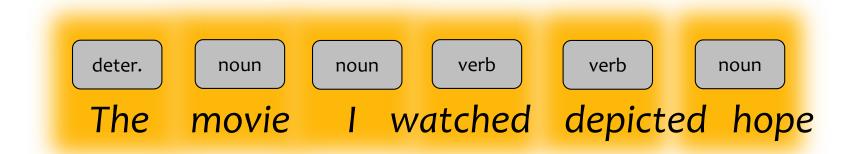


Table from Manning (2011)

#### Feature Engineering for NLP

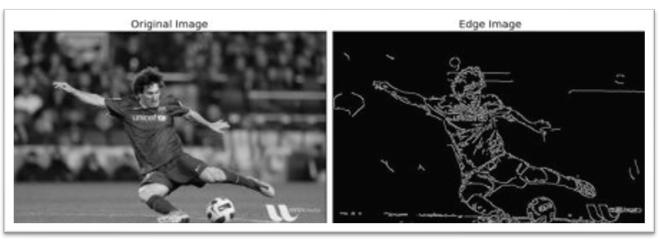
**Table 3.** Tagging accuracies with different feature templates and other changes on the WSJ 19-21 development set.

			7		
Model	Feature Templates	#	Sent.	Token	Unk.
		Feats	Acc.	Acc.	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%
$5\mathrm{w}$	$+\langle t_0,w_{-2} angle, \langle t_0,w_2 angle$	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%

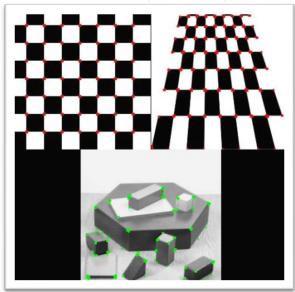
deter. noun verb verb noun

The movie I watched depicted hope

#### Edge detection (Canny)



#### Corner Detection (Harris)



#### Scale Invariant Feature Transform (SIFT)

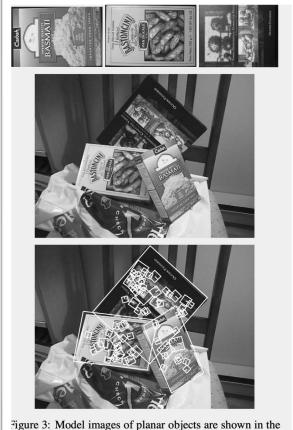


Figure 3: Model images of planar objects are shown in the oprow. Recognition results below show model outlines and mage keys used for matching.

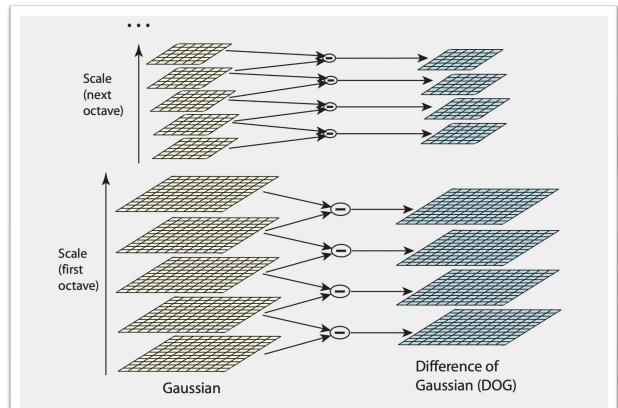


Figure 1: For each octave of scale space, the initial image is repeatedly convolved with Gaussians to produce the set of scale space images shown on the left. Adjacent Gaussian images are subtracted to produce the difference-of-Gaussian images on the right. After each octave, the Gaussian image is down-sampled by a factor of 2, and the process repeated.

#### **NON-LINEAR FEATURES**

#### Nonlinear Features

- aka. "nonlinear basis functions"
- So far, input was always  $\mathbf{x} = [x_1, \dots, x_M]$
- **Key Idea:** let input be some function of **x** 
  - $\begin{array}{ll} & \text{original input:} & \mathbf{x} \in \mathbb{R}^M \\ & \text{new input:} & \mathbf{x}' \in \mathbb{R}^{M'} \end{array} \text{ where } M' > M \text{ (usually)}$

  - define  $\mathbf{x}' = b(\mathbf{x}) = [b_1(\mathbf{x}), b_2(\mathbf{x}), \dots, b_{M'}(\mathbf{x})]$

where  $b_i: \mathbb{R}^M \to \mathbb{R}$  is any function

Examples: (M = 1)

$$b_j(x) = x^j \quad \forall j \in \{1, \dots, J\}$$

radial basis function

$$b_j(x) = \exp\left(\frac{-(x-\mu_j)^2}{2\sigma_j^2}\right)$$

sigmoid

$$b_j(x) = \frac{1}{1 + \exp(-\omega_j x)}$$

log

$$b_j(x) = \log(x)$$

#### For a linear model: still a linear function of b(x) even though a

nonlinear function of

X

#### **Examples:**

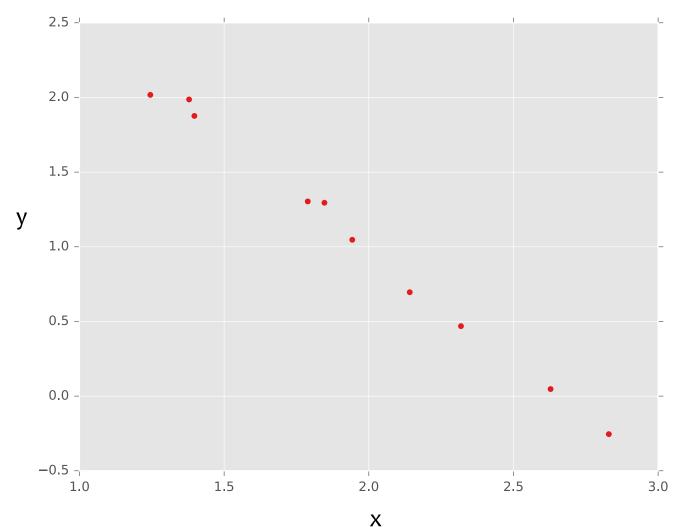
- Perceptron
- Linear regression
- Logistic regression

#### Example: Linear Regression

**Goal:** Learn  $y = \mathbf{w}^T f(\mathbf{x}) + \mathbf{b}$  where f(.) is a polynomial basis function

i	у	х
1	2.0	1.2
2	1.3	1.7
•••	•••	•••
10	1.1	1.9

true "unknown"
target function is
linear with
negative slope
and gaussian
noise

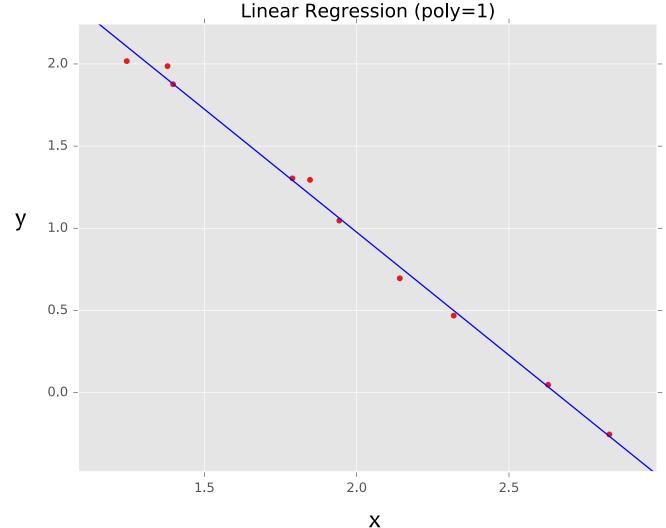


#### Example: Linear Regression

**Goal:** Learn  $y = \mathbf{w}^T f(\mathbf{x}) + b$  where f(.) is a polynomial basis function

i	у	х
1	2.0	1.2
2	1.3	1.7
•••		•••
10	1.1	1.9

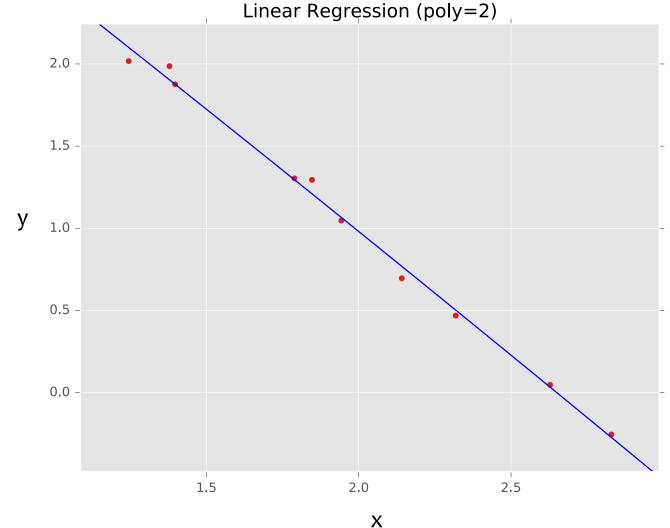
true "unknown"
target function is
linear with
negative slope
and gaussian
noise



**Goal:** Learn  $y = \mathbf{w}^T f(\mathbf{x}) + b$  where f(.) is a polynomial basis function

i	у	х	X <sup>2</sup>
1	2.0	1.2	(1.2)2
2	1.3	1.7	(1.7)2
•••	•••	•••	•••
10	1.1	1.9	(1.9)2

true "unknown"
target function is
linear with
negative slope
and gaussian
noise

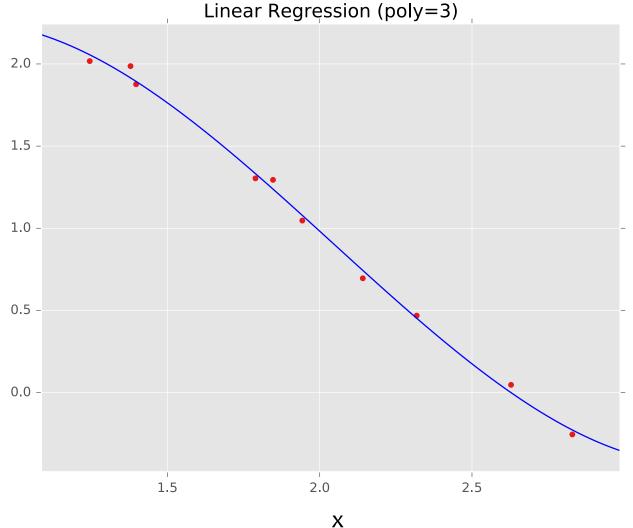


**Goal:** Learn  $y = \mathbf{w}^T f(\mathbf{x}) + \mathbf{b}$  where f(.) is a polynomial

basis function

i	у	х	X <sup>2</sup>	<b>X</b> <sup>3</sup>
1	2.0	1.2	(1.2)2	(1.2)3
2	1.3	1.7	(1.7)2	<b>(1.7)</b> <sup>3</sup>
•••	•••	•••	•••	•••
10	1.1	1.9	(1.9)2	(1.9)3

true "unknown"
target function is
linear with
negative slope
and gaussian
noise

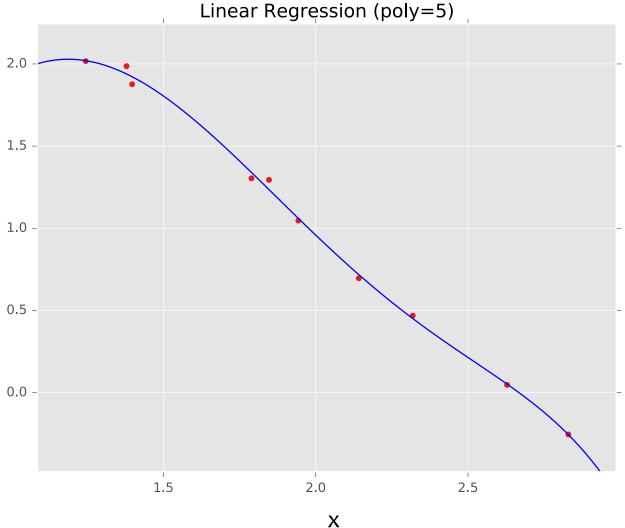


**Goal:** Learn  $y = \mathbf{w}^T f(\mathbf{x}) + \mathbf{b}$  where f(.) is a polynomial

basis function

i	у	х		<b>x</b> <sup>5</sup>
1	2.0	1.2	•••	(1.2)5
2	1.3	1.7	•••	(1.7)5
•••		•••		•••
10	1.1	1.9	•••	(1.9)5

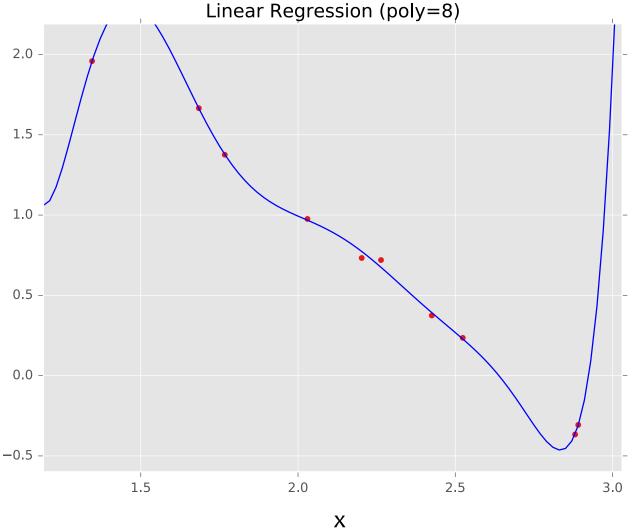
true "unknown"
target function is
linear with
negative slope
and gaussian
noise



**Goal:** Learn  $y = \mathbf{w}^T f(\mathbf{x}) + b$  where f(.) is a polynomial basis function

i	у	х	•••	X <sup>8</sup>
1	2.0	1.2		(1.2)8
2	1.3	1.7	•••	(1.7)8
	•••	•••		•••
10	1.1	1.9	•••	(1.9)8

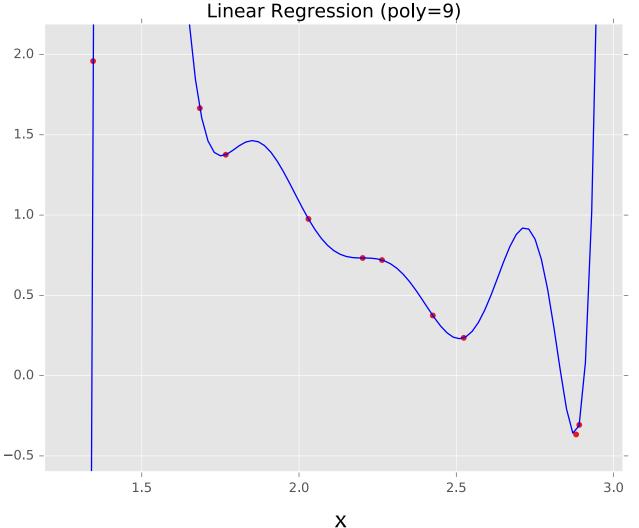
true "unknown" target function is linear with negative slope and gaussian noise



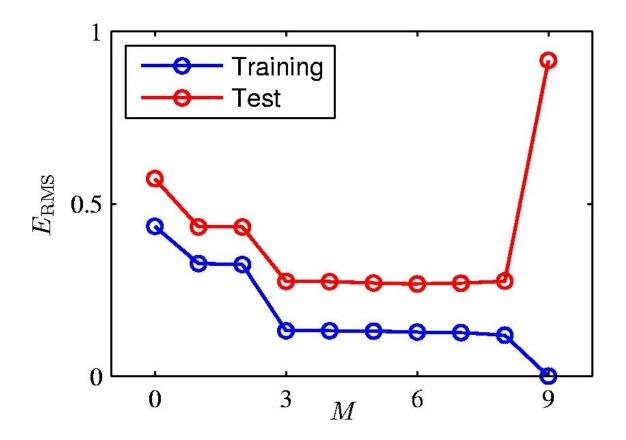
**Goal:** Learn  $y = \mathbf{w}^T f(\mathbf{x}) + b$  where f(.) is a polynomial basis function

i	у	х		<b>x</b> <sup>9</sup>
1	2.0	1.2	•••	(1.2)9
2	1.3	1.7	•••	(1.7)9
	•••	•••	•••	•••
10	1.1	1.9	•••	(1.9)9

true "unknown" target function is linear with negative slope and gaussian noise



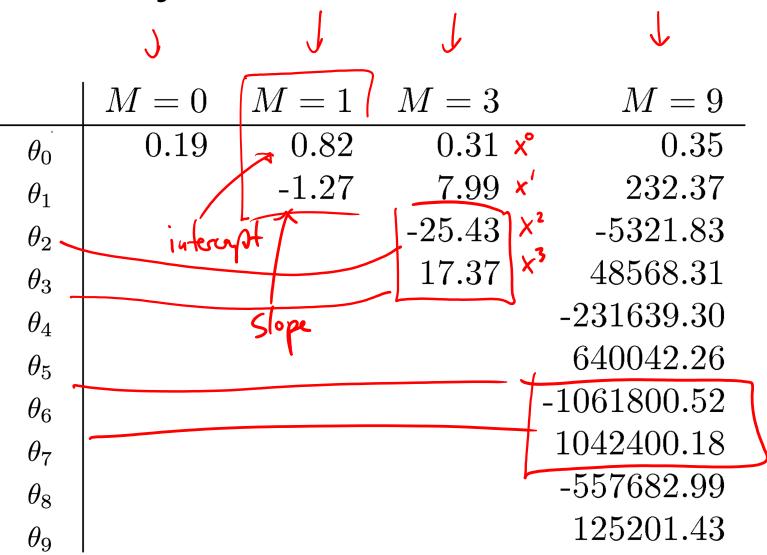
# Over-fitting



Root-Mean-Square (RMS) Error:  $E_{\rm RMS} = \sqrt{2E(\mathbf{w}^{\star})/N}$ 

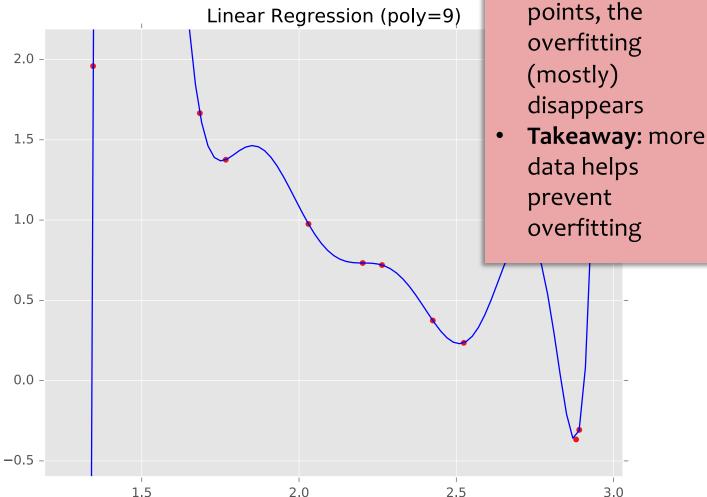
$$E_{\rm RMS} = \sqrt{2E(\mathbf{w}^{\star})/N}$$

## Polynomial Coefficients



**Goal:** Learn  $y = \mathbf{w}^T f(\mathbf{x}) + b$  where f(.) is a polynomial basis function

i	у	х	•••	<b>x</b> <sup>9</sup>	
1	2.0	1.2	•••	(1.2)9	
2	1.3	1.7	•••	(1.7)9	
	•••	•••	•••		у
10	1.1	1.9	•••	(1.9)9	



X

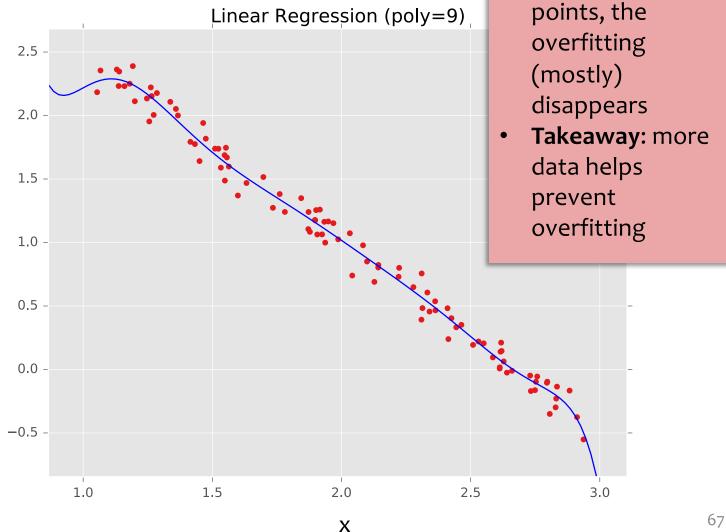
With just N = 10

But with N = 100

points we overfit!

**Goal:** Learn  $y = \mathbf{w}^T f(\mathbf{x}) + b$  where f(.) is a polynomial basis function

i	у	х	•••	<b>x</b> <sup>9</sup>	
1	2.0	1.2	•••	(1.2)9	
2	1.3	1.7	•••	<b>(1.7)</b> <sup>9</sup>	
3	0.1	2.7		(2.7)9	,
4	1.1	1.9	•••	(1.9)9	
	•••	•••			
	•••	•••	•••		
		•••			
98	•••	•••			
99	•••	•••	•••	•••	
100	0.9	1.5	•••	(1.5)9	



With just N = 10

But with N = 100

points we overfit!

#### **REGULARIZATION**

# Overfitting

**Definition:** The problem of **overfitting** is when the model captures the noise in the training data instead of the underlying structure

Overfitting can occur in all the models we've seen so far:

- Decision Trees (e.g. when tree is too deep)
- KNN (e.g. when k is small)
- Perceptron (e.g. when sample isn't representative)
- Linear Regression (e.g. with nonlinear features)
- Logistic Regression (e.g. with many rare features)

## Motivation: Regularization

 Occam's Razor: prefer the simplest hypothesis

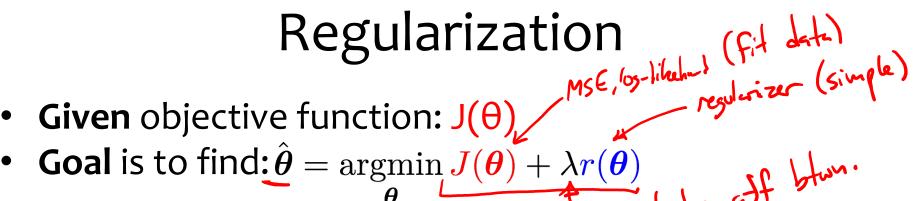
- What does it mean for a hypothesis (or model) to be simple?
  - 1. small number of features (model selection)
  - small number of "important" features (shrinkage)

$$\vec{\Theta} = \begin{bmatrix} \Theta_1 \\ \Theta_2 \\ \Theta_3 \\ \Theta_4 \end{bmatrix} = \begin{bmatrix} 23 \\ 7 \\ 0.0000001 \\ -4 \end{bmatrix}$$

$$\vec{X} = \begin{bmatrix} 1 \\ 2 \\ -1 \\ 4 \end{bmatrix}$$

$$\vec{\Theta} \vec{X} = \begin{bmatrix} 1 \\ 2 \\ -1 \\ 4 \end{bmatrix}$$

$$\vec{\Theta} \vec{X} = \begin{bmatrix} 1 \\ 2 \\ -1 \\ 4 \end{bmatrix}$$



• 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  $r(\theta)$ :

Example: q-norm (usually p-norm):  $\|\boldsymbol{\theta}\|_q = \left(\sum_{m=1}^{M} |\theta_m|^q\right)^{\overline{q}}$ 

(A)	q -2	$r(oldsymbol{ heta})$	yields parame- ters that are	name	optimization notes
9	0	$  \boldsymbol{\theta}  _0 = \sum \mathbb{1}(\theta_m \neq 0)$	zero values	Lo reg.	no good computa- tional solutions
_	1 2	$  oldsymbol{ heta}  _1 = \sum   heta_m  \ (  oldsymbol{ heta}  _2)^2 = \sum  heta_m^2$	zero values small values	L1 reg. L2 reg.	subdifferentiable differentiable