



10-301/10-601 Introduction to Machine Learning

Machine Learning Department School of Computer Science Carnegie Mellon University

Logistic Regression + Feature Engineering + Regularization

Matt Gormley & Henry Chai Lecture 10 Feb. 12, 2025

Reminders

- Practice Problems 1
 - released on course website
- Exam 1: Mon, Feb. 17
 - Time: 7:00 9:00pm
 - Location: Your room/seat assignment will be announced on Piazza
- Homework 4: Logistic Regression
 - Out: Mon, Feb 17
 - Due: Wed, Feb. 26 at 11:59pm

EXAM 1 LOGISTICS

Exam 1

Time / Location

- Time: Mon, Feb 17, at 7:00pm 9:00pm
- Location & Seats: You have all been split across multiple rooms. Everyone has an assigned seat in one of these room.
- Please watch Piazza carefully for announcements.

Logistics

- Covered material: Lecture 1 Lecture 7
- Format of questions:
 - Multiple choice
 - True / False (with justification)
 - Derivations
 - Short answers
 - Interpreting figures
 - Implementing algorithms on paper
- No electronic devices
- You are allowed to **bring** one 8% x 11 sheet of notes (front and back)

Exam 1

How to Prepare

- Attend the Exam OHs on Friday
- Review exam practice problems
- Review this year's homework problems
- Consider whether you have achieved the "learning objectives" for each lecture / section
- Write your one-page cheat sheet (back and front)

Exam 1

Advice (for during the exam)

- Solve the easy problems first
 (e.g. multiple choice before derivations)
 - if a problem seems extremely complicated you're likely missing something
- Don't leave any answer blank!
- If you make an assumption, write it down
- If you look at a question and don't know the answer:
 - we probably haven't told you the answer
 - but we've told you enough to work it out
 - imagine arguing for some answer and see if you like it

Topics for Exam 1

- Foundations
 - Probability, Linear Algebra,
 Geometry, Calculus
 - Optimization
- Important Concepts
 - Overfitting
 - Experimental Design

- Classification
 - Decision Tree
 - KNN
 - Perceptron
- Regression
 - KNN Regression
 - Decision Tree Regression
 - Linear Regression

LOGISTIC REGRESSION

1. Model

$$y \sim \text{Bernoulli}(\emptyset)$$

$$\emptyset = \sigma(\widehat{\Theta}^{T}\widehat{x}) \text{ where } \sigma(\upsilon) = \frac{1}{1+\exp(-\upsilon)}$$

$$p(y|\widehat{x},\Theta) = \int \sigma(\widehat{\Theta}^{T}\widehat{x}) : f(y=0)$$

$$(1-\sigma(\widehat{\Theta}^{T}\widehat{x})) : f(y=0)$$

2. Objective

$$l(\Theta) = \log p(D|\overline{\Theta}) = \log \prod_{i=1}^{N} p(y^{(i)}|\dot{x}^{(i)}, \overline{\Theta})$$

$$= \sum_{i=1}^{N} \log p(y^{(i)}|\dot{x}^{(i)}, \overline{\Theta})$$

$$J(\theta) = -\frac{1}{N} \mathcal{N}(\vec{\theta})$$

$$= \frac{1}{N} \sum_{i=1}^{N} -\log p(y^{(i)} | \vec{x}^{(i)}, \vec{\theta})$$

$$J(\theta) = -\frac{1}{N} \mathcal{N}(\vec{\theta})$$

3A. Derivatives

3B. Gradients

$$\nabla \mathcal{J}^{(i)}(\theta) = \left[\begin{array}{c} -(\gamma^{(i)} - \sigma(\widehat{\theta}^{T}\widehat{x}^{(i)})) \\ \vdots \\ \vdots \\ \end{array} \right]$$

4. Optimization

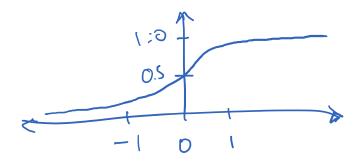
5. Prediction

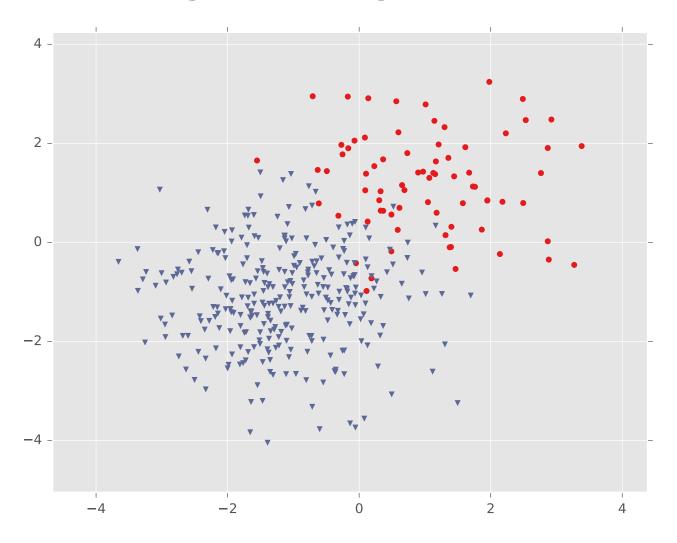
Predict
$$\hat{y} = \underset{\text{argmax}}{\text{argmax}} P(y|\hat{x})$$

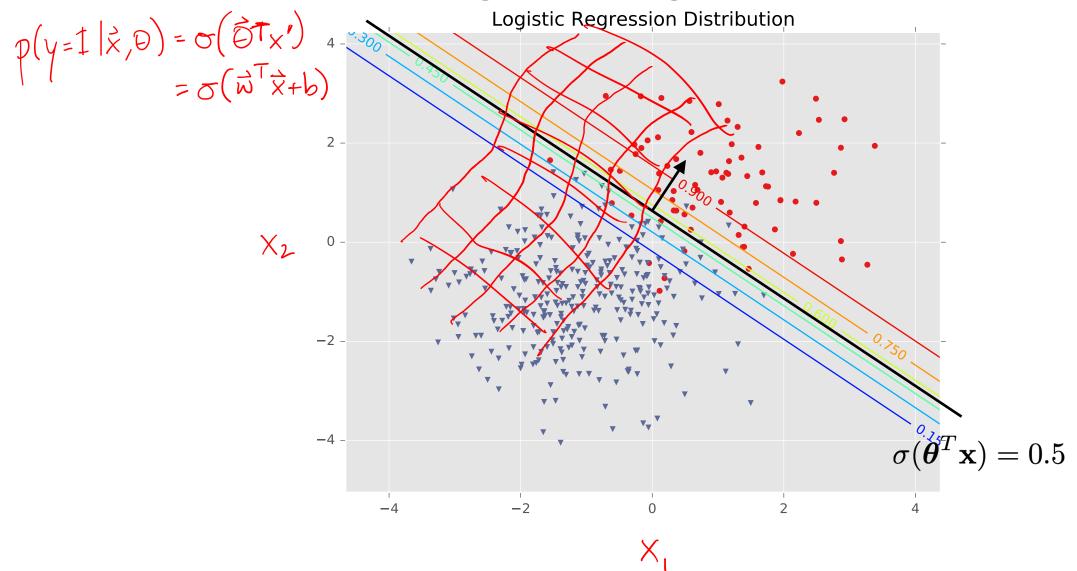
$$= \cdots$$

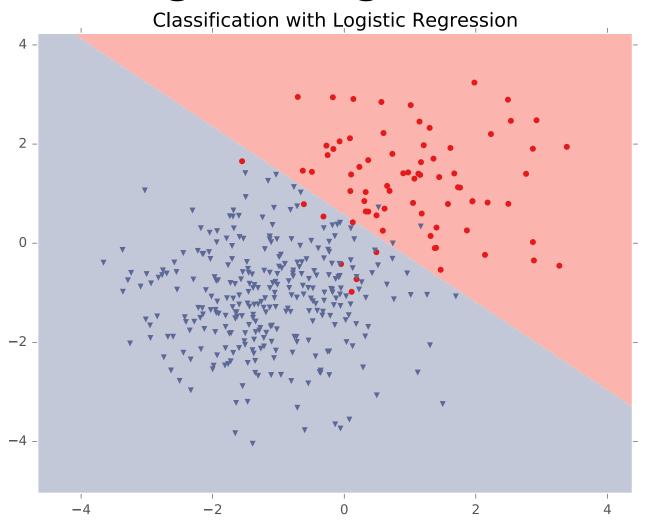
$$= \cdots$$

$$= "Sign"(\hat{\Theta}^{T}\hat{x})$$









Example: Image Classification

- ImageNet LSVRC-2010 contest:
 - Dataset: 1.2 million labeled images, 1000 classes
 - Task: Given a new image, label it with the correct class
 - Multiclass classification problem
- Examples from http://image-net.org/

Bird

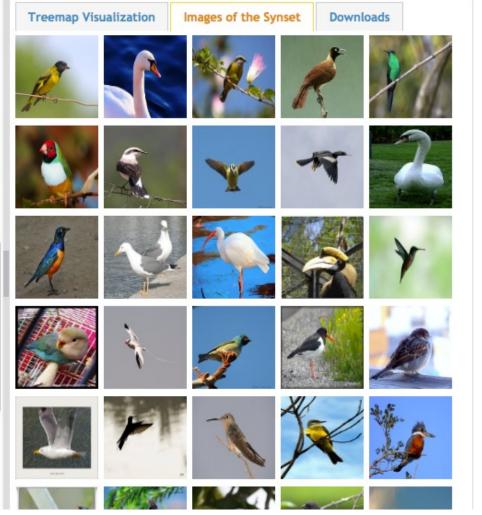
IM ... GENET

Warm-blooded egg-laying vertebrates characterized by feathers and forelimbs modified as wings

2126 pictures 92.85% Popularity Percentile



marine animal, marine creature, sea animal, sea creature (1)
⊩ scavenger (1)
- biped (0)
redator, predatory animal (1)
⊩ larva (49)
- acrodont (0)
- feeder (0)
stunt (0)
r-chordate (3087)
tunicate, urochordate, urochord (6)
- cephalochordate (1)
vertebrate, craniate (3077) vertebrate, craniate (3077)
mammal, mammalian (1169)
∳- bird (871)
dickeybird, dickey-bird, dickybird, dicky-bird (0)
⊩- cock (1)
- hen (0)
nester (0)
night bird (1)
- bird of passage (0)
- protoavis (0)
archaeopteryx, archeopteryx, Archaeopteryx lithographi
- Sinornis (0)
- Ibero-mesornis (0)
- archaeornis (0)
ratite, ratite bird, flightless bird (10)
carinate, carinate bird, flying bird (0)
passerine, passeriform bird (279)
- nonpasserine bird (0)
bird of prey, raptor, raptorial bird (80)
gallinaceous bird, gallinacean (114)



Not logged in. Login I Signup

German iris, Iris kochii

Iris of northern Italy having deep blue-purple flowers; similar to but smaller than Iris germanica

469 pictures 49.6% Popularity Percentile







Not logged in. Login I Signup

Court, courtyard

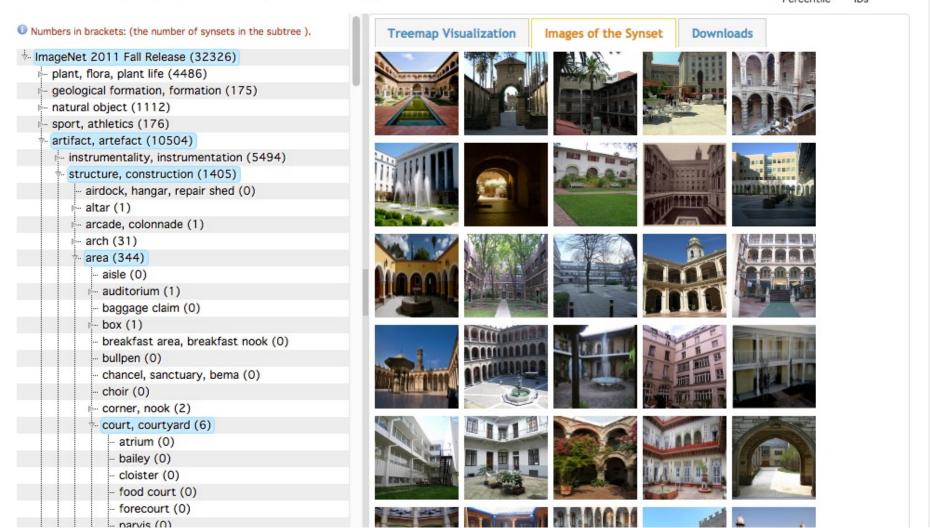
IM GENET

An area wholly or partly surrounded by walls or buildings; "the house was built around an inner court"

165 pictures

92.61% Popularity Percentile





Example: Image Classification

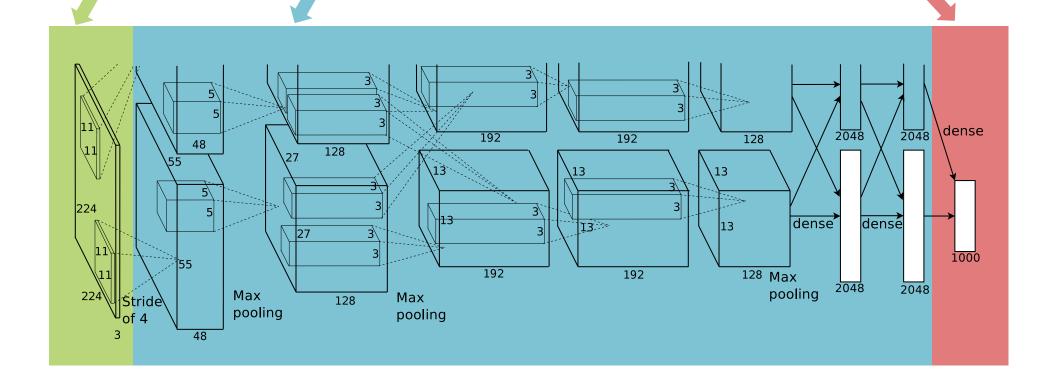
CNN for Image Classification

(Krizhevsky, Sutskever & Hinton, 2011) 17.5% error on ImageNet LSVRC-2010 contest

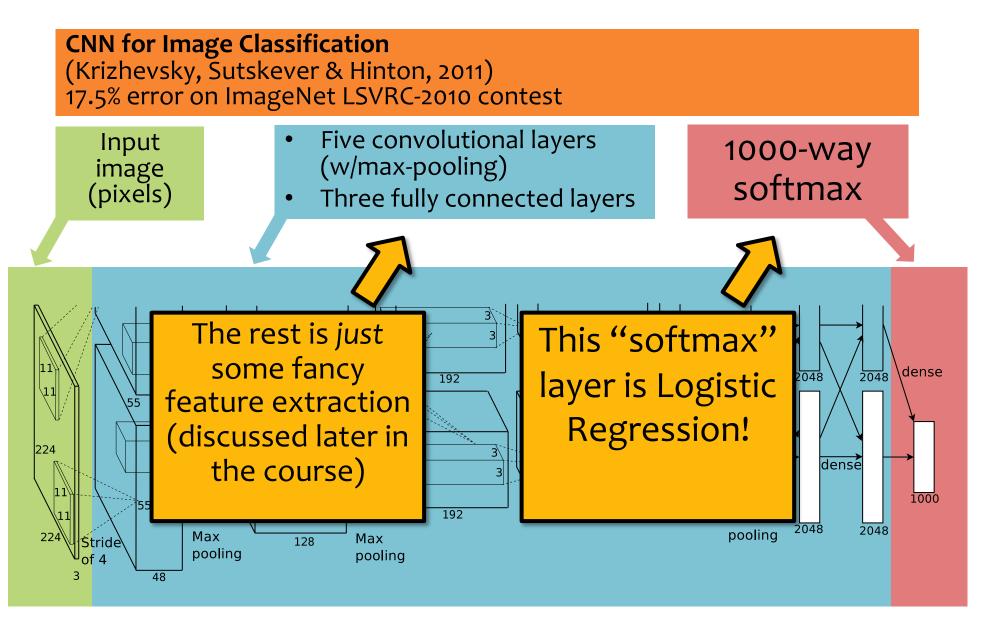
Input image (pixels)

- Five convolutional layers (w/max-pooling)
- Three fully connected layers

1000-way softmax



Example: Image Classification



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 classification
- Prove that the decision boundary of binary logistic regression is linear

Linear Models

PERCEPTRON, LINEAR REGRESSION, AND LOGISTIC REGRESSION

Matching Game uplate for one porceder

Poll Question:

Match the Algorithm to its Update Rule

1. SGD for Logistic Regression

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = p(y = 1 \mid \mathbf{x})$$

2. Least Mean Squares

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

3. Perceptron

$$\frac{h_{\boldsymbol{\theta}}(\mathbf{x}) = \operatorname{sign}(\boldsymbol{\theta}^T \mathbf{x})}{(\mathbf{x})^2}$$

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

5.
$$\theta_k \leftarrow \theta_k + \frac{1}{1 + \exp \lambda (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)})}$$

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

Answer:

SGD for

Linew

I. None of the above



Gradient Descent

Algorithm 1 Gradient Descent

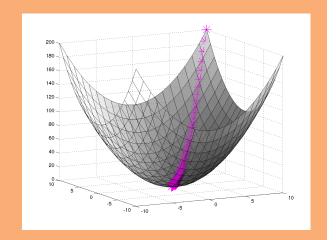
1: **procedure** $GD(\mathcal{D}, \boldsymbol{\theta}^{(0)})$

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 θ



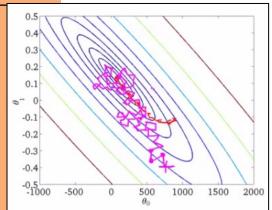
In order to apply GD to Logistic Regression all we need is the **gradient** of the objective function (i.e. vector of partial derivatives).

$$egin{align} egin{align} egin{align} rac{d}{d heta_1}J(oldsymbol{ heta}) \ rac{d}{d heta_2}J(oldsymbol{ heta}) \ dots \ rac{d}{d heta_M}J(oldsymbol{ heta}) \ \end{pmatrix} \end{aligned}$$

Stochastic Gradient Descent (SGD)

Algorithm 1 Stochastic Gradient Descent (SGD)

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



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

We need a per-example objective:

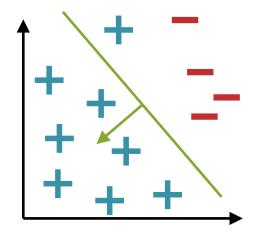
Let
$$J(\boldsymbol{\theta}) = \sum_{i=1}^N J^{(i)}(\boldsymbol{\theta})$$
 where $J^{(i)}(\boldsymbol{\theta}) = -\log p_{\boldsymbol{\theta}}(y^i|\mathbf{x}^i)$.

Logistic Regression vs. Perceptron

Poll 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.

Answer:



BAYES OPTIMAL CLASSIFIER

Bayes Optimal Classifier

was gen function

Suppose you knew the Functio 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?

Answer:

Our goa You'd use the Bayes best app optimal classifier!

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 $p^*(y|x)$

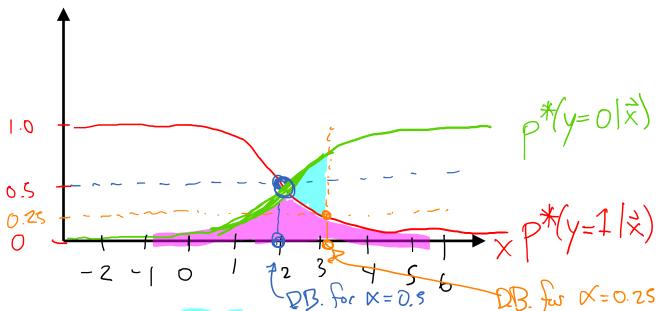
Bayes Optimal Classifier

Suppose you have an **oracle** that knows the data generating distribution, p*(y|x).

Q: What is the optimal classifier in this setting?

A: The Bayes optimal classifier! This is the best classifier for the distribution p* and

the loss function.



Definition: The **reducible error** is the expected loss of a hypothesis h(x) that could be reduced if we knew p*(y|x) and picked the optimal h(x) for that p*.

Definition: The irreducible error is the expected loss of a hypothesis h(x) that could **not** be reduced if we knew p*(y|x) and picked the optimal h(x) for that p*.

$$J(y,\hat{y}) = I(y \neq \hat{y})$$

$$\hat{Y} = h(x) = \begin{cases} 1 & \text{if } \hat{y}(y \neq \hat{y}) \geq x \\ 0 & \text{otherwise} \end{cases}$$

$$x = 0.5 \text{ for } 0/1 \text{ loss}$$

$$|= \begin{cases} 1 \text{ million if } y \neq \hat{y} \text{ and } y = 1 \\ 1000 \text{ if } y \neq \hat{y} \text{ and } y = 0 \\ 0 \text{ otherwise} \end{cases}$$

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 **%** 5 randomly chosen examples

Mini-Batch SGD

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

Three variants of first-order optimization:

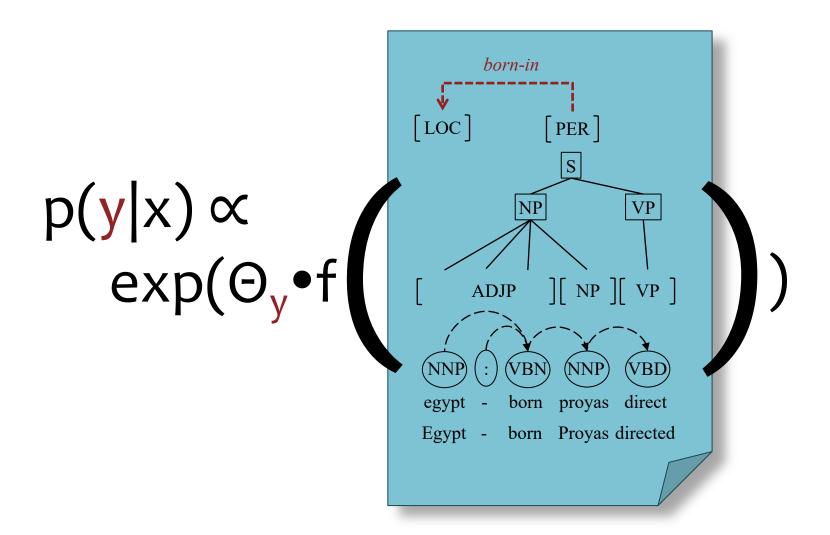
Gradient Descent:
$$\mathbf{g} = \nabla J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^{N} \nabla J^{(i)}(\boldsymbol{\theta})$$

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

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

FEATURE ENGINEERING

Handcrafted Features



Where do features come from?

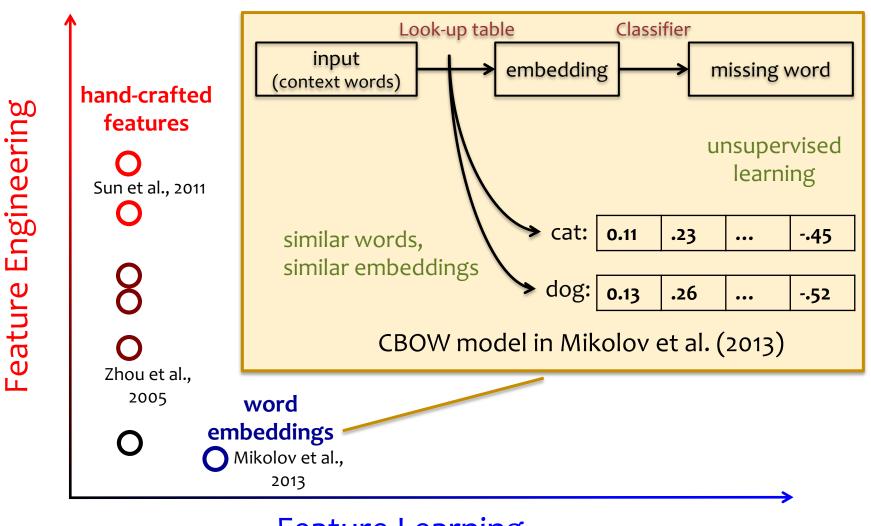
First word before M1 Second word before M1 hand-crafted Bag-of-words in M1 features Head word of M1 Other word in between *First word after M2* Sun et al., 2011 Second word after M2 Bag-of-words in M2 *Head word of M2* Bigrams in between Words on dependency path Country name list Personal relative triggers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

Engineering

Feature

Feature Learning

Where do features come from?

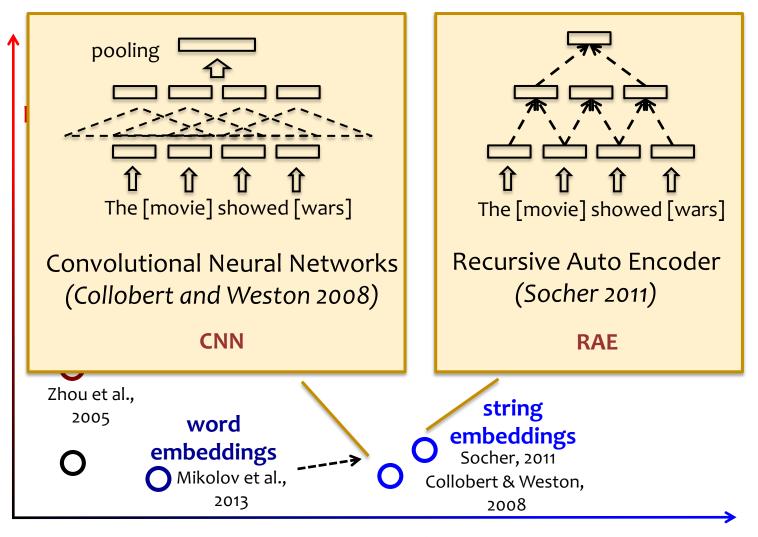


Feature Learning

Engineering

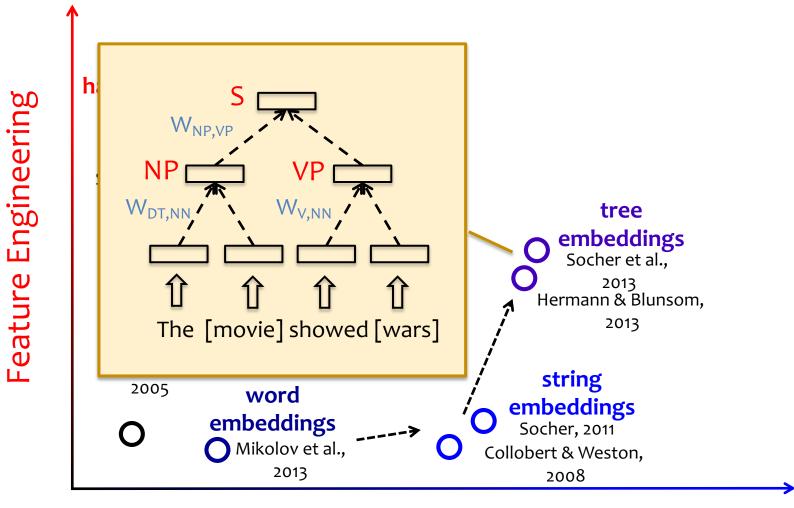
Feature

Where do features come from?



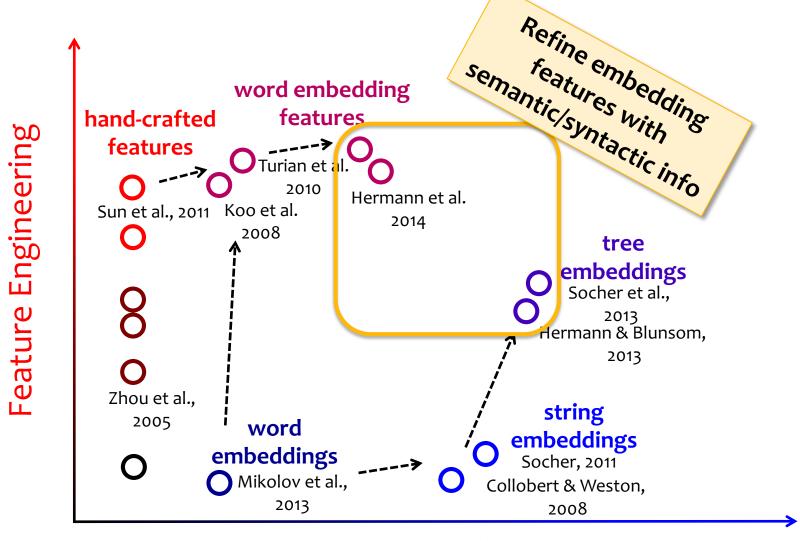
Feature Learning

Where do features come from?



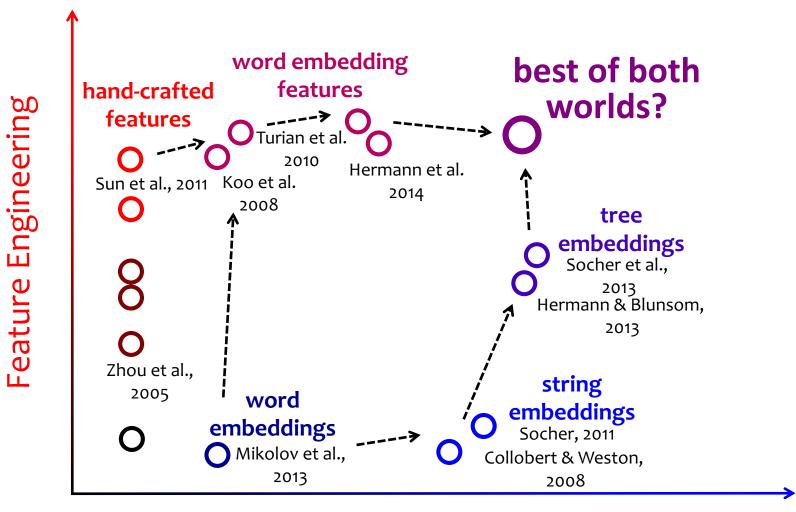
Feature Learning

Where do features come from?



Feature Learning

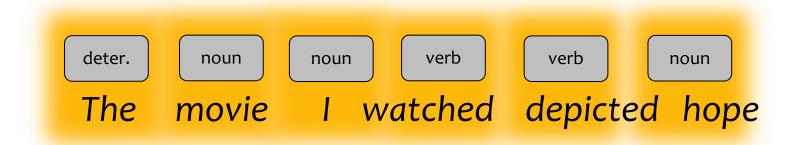
Where do features come from?



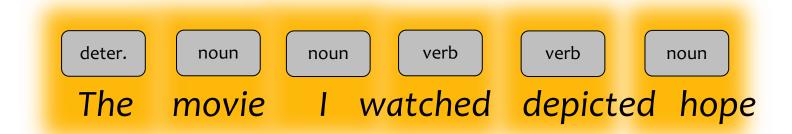
Feature Learning

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

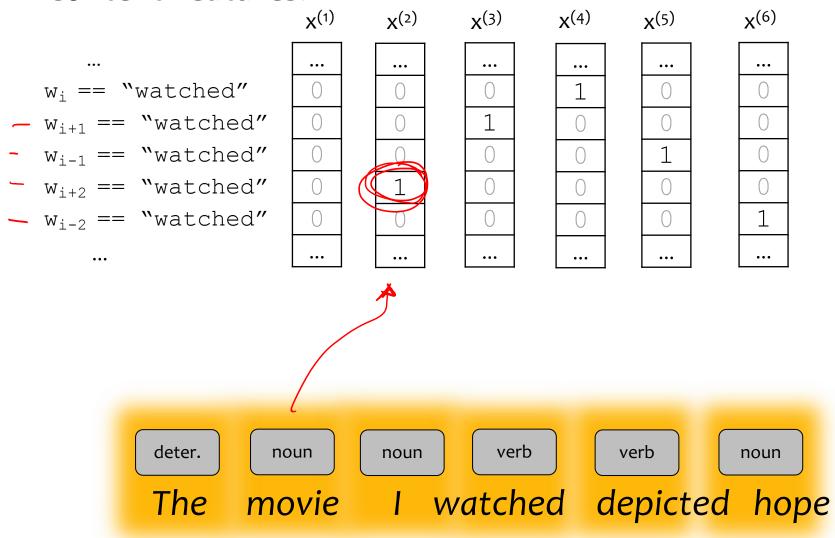
What features should you use?



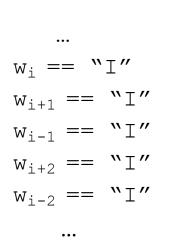
Per-word Features:



Context Features:



Context Features:



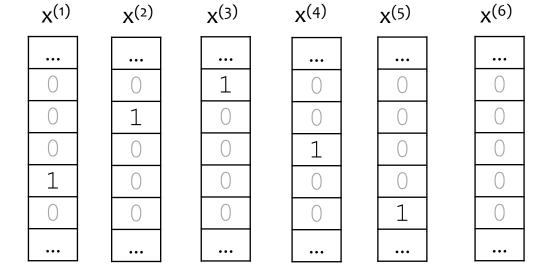




Table 3. Tagging accuracies with different feature templates and other changes on the WSJ 19-21 development set.

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%

The movie I watched depicted hope

Background: Word Embeddings

One-hot vectors

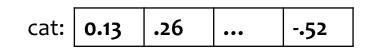
- Standard representation of a word in NLP:
 1-hot vector (aka. a string)
- Vectors representing related words share nothing in common

	٥	and	pe	ω×	90go		401	zebra
cat:	0	0	0	1	0	•••	0	0
dog:	0	0	0	0	1	•••	0	0

Word embeddings

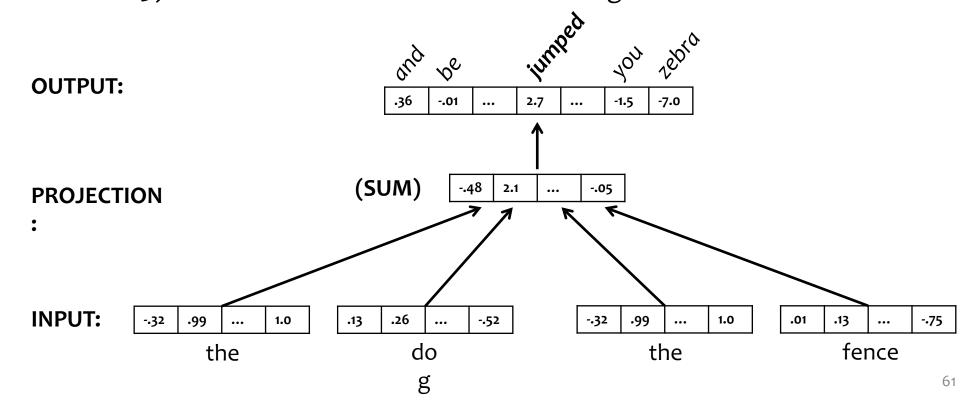
M = 1024

- 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



Background: Word Embeddings

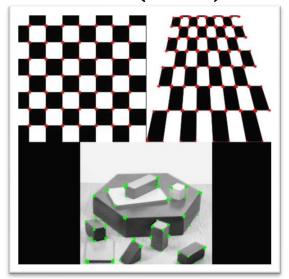
- It's common to use neural-network trained embeddings
 - Key idea: learn embeddings which are good at reconstructing the context of a word
 - Popular across HLT (speech, NLP)
- The Continuous Bag-of-words Model (CBOW) (Mikolov et al., 2013) maximizes the likelihood of a word given its context:



Edge detection (Canny)



Corner Detection (Harris)



Figures from http://opencv.org

Scale Invariant Feature Transform (SIFT)



mage keys used for matching.

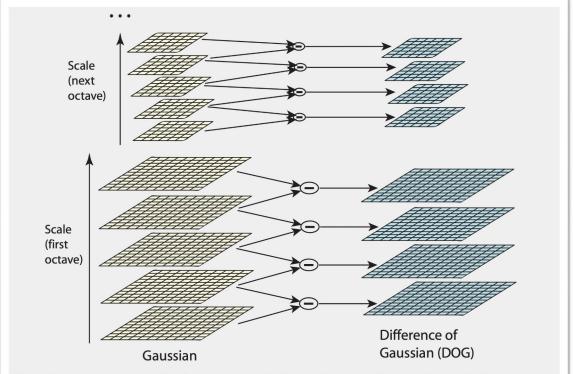


Figure 1: For each octave of scale space, the initial image is repeatedly convolved with Gaussians to produce the set of scale space images shown on the left. Adjacent Gaussian images are subtracted to produce the difference-of-Gaussian images on the right. After each octave, the Gaussian image is down-sampled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

NON-LINEAR FEATURES

Nonlinear Features

- aka. "nonlinear basis functions"
- So far, input was always $\mathbf{x} = [x_1, \dots, x_M]$

$$\mathbf{x} = [x_1, \dots, x_M]$$

- **Key Idea:** let input be some function of **x**
 - original input:

$$\mathbf{x} \in \mathbb{R}^M$$

 $\mathbf{x} \in \mathbb{R}^M$ where M' > M (usually) $\mathbf{x}' \in \mathbb{R}^{M'}$

– new input:

$$\mathbf{x}' \in \mathbb{R}^{M'}$$

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

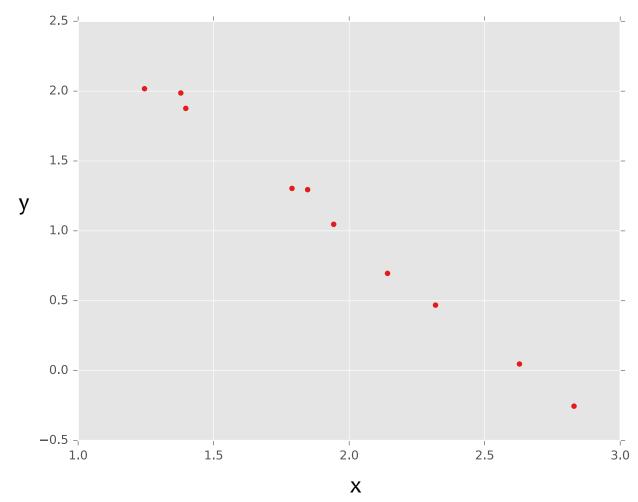
X

Examples:

- Perceptron
- Linear regression
- Logistic regression

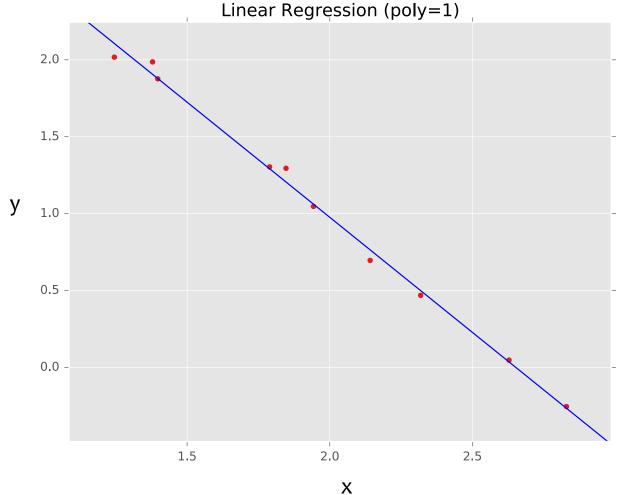
Goal: Learn $y = \mathbf{w}^T f(\mathbf{x}) + b$ where f(.) is a polynomial basis function

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



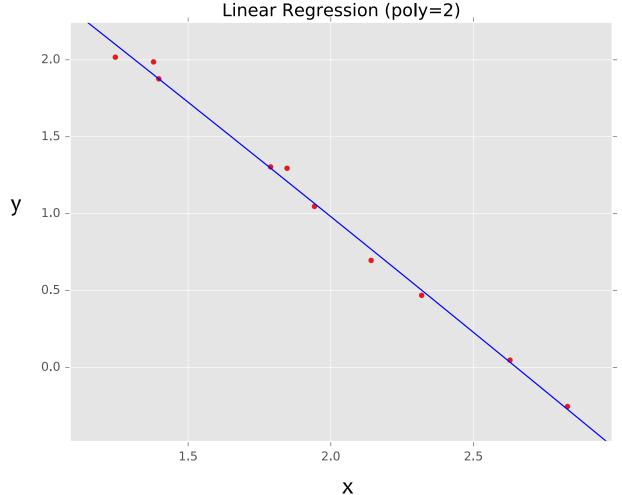
Goal: Learn $y = \mathbf{w}^T f(\mathbf{x}) + b$ where f(.) is a polynomial basis function

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



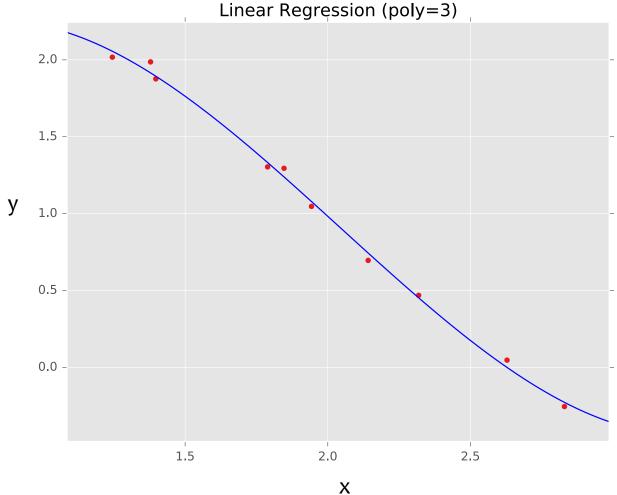
Goal: Learn $y = \mathbf{w}^T f(\mathbf{x}) + b$ where f(.) is a polynomial basis function

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



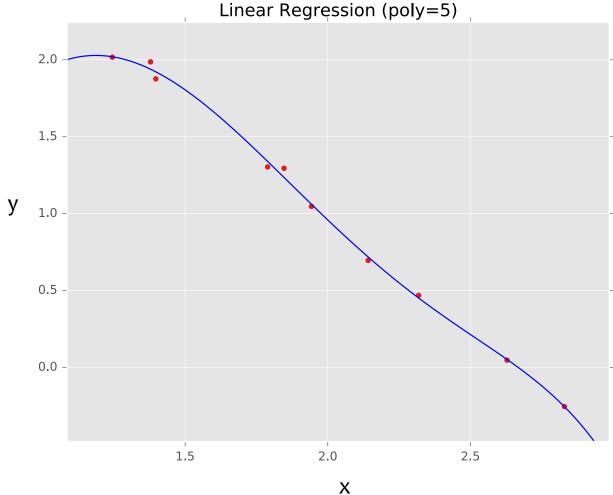
Goal: Learn $y = \mathbf{w}^T f(\mathbf{x}) + b$ where f(.) is a polynomial basis function

i	у	х	X ²	x ³
1	2.0	1.2	(1.2)2	(1.2)3
2	1.3	1.7	(1.7)2	(1.7)3
•••	•••		•••	•••
10	1.1	1.9	(1.9)2	(1.9)3



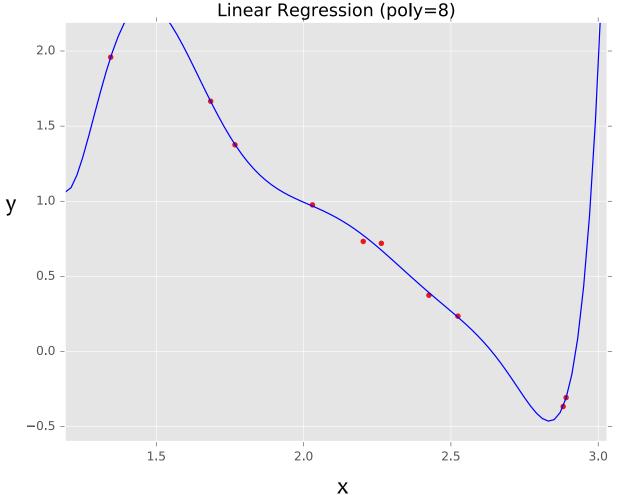
Goal: Learn $y = \mathbf{w}^T f(\mathbf{x}) + b$ where f(.) is a polynomial basis function

i	у	х	•••	x ⁵
1	2.0	1.2	•••	(1.2)5
2	1.3	1.7	•••	(1.7)5
		•••	•••	•••
10	1.1	1.9	•••	(1.9)5



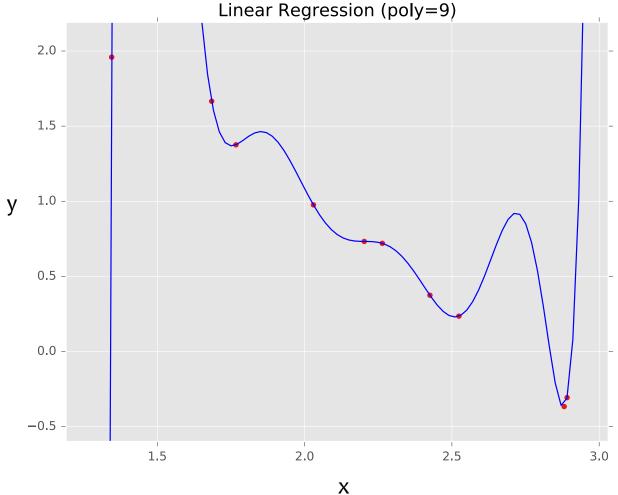
Goal: Learn $y = \mathbf{w}^T f(\mathbf{x}) + b$ where f(.) is a polynomial basis function

i	у	х	•••	x ⁸
1	2.0	1.2	•••	(1.2)8
2	1.3	1.7		(1.7)8
	•••	•••		•••
10	1.1	1.9		(1.9)8

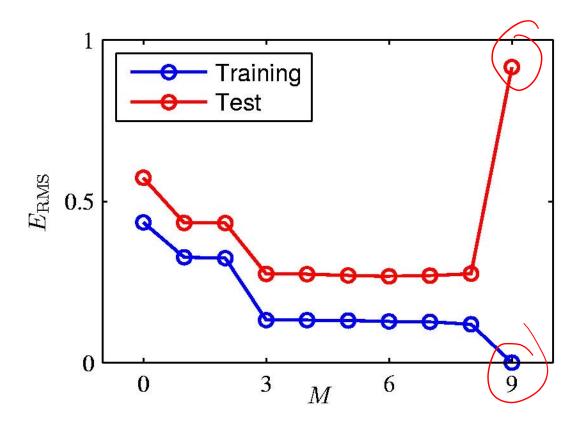


Goal: Learn $y = \mathbf{w}^T f(\mathbf{x}) + b$ where f(.) is a polynomial basis function

i	у	х	•••	x ⁹
1	2.0	1.2	•••	(1.2)9
2	1.3	1.7		(1.7)9
•••	•••	•••		
10	1.1	1.9	•••	(1.9)9

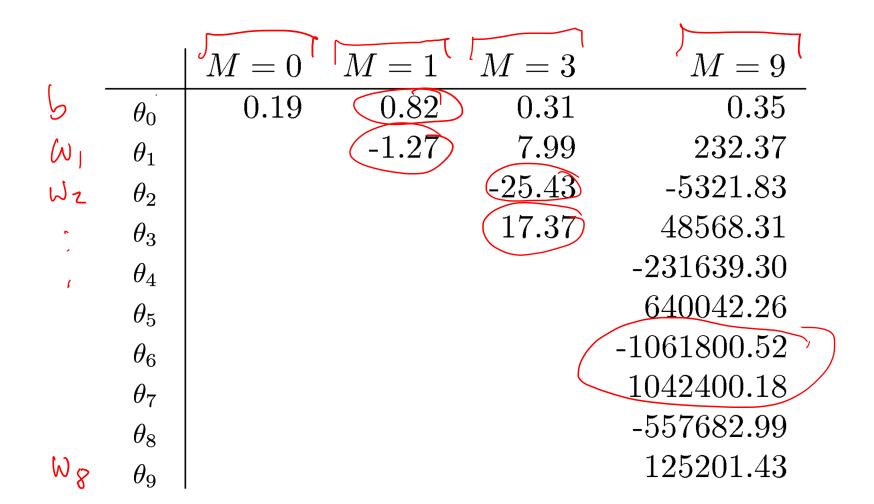


Over-fitting



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

Polynomial Coefficients



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

2.0

1.5 -

1.0 -

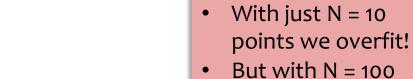
0.5 -

0.0 -

-0.5 -

1.5

i	у	х		x ⁹	
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	



- But with N = 1
 points, the overfitting (mostly) disappears
- Takeaway: more data helps prevent overfitting

3.0

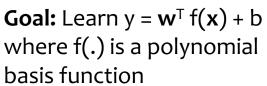
2.5

2.0

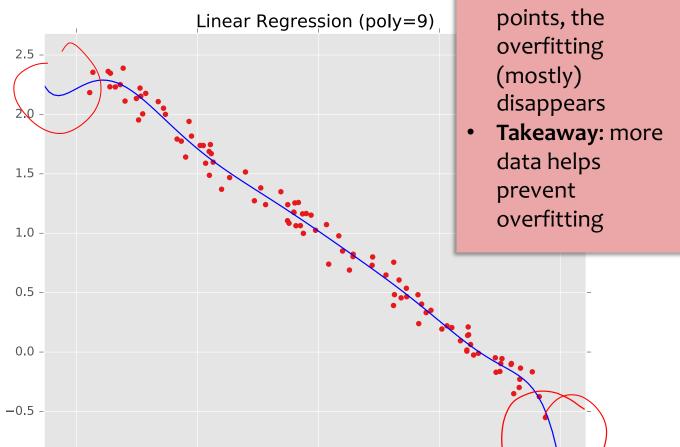
Linear Regression (poly=9)

1.5

1.0



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



2.0

Χ

2.5

• With just N = 10

points we overfit!

But with N = 100

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)
 - 2. small number of "important" features (shrinkage)

$$\overrightarrow{X} = \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \\ X_7 \end{bmatrix}$$

$$\overrightarrow{X} = \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \\ X_5 \\ X_6 \\ X_7 \end{bmatrix}$$

$$\overrightarrow{X} = \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \\ X_5 \\ X_6 \\ X_7 \end{bmatrix}$$

Regularization

- **Given** objective function: $J(\theta)$
- **Goal** is to find:

$$\hat{\boldsymbol{\theta}} = \operatorname*{argmin}_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) + \lambda r(\boldsymbol{\theta})$$

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

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