# Logistic Regression + Feature Engineering + Regularization

Matt Gormley & Henry Chai
Lecture 10
Oct. 1, 2021

# Reminders

- **Homework 4: Logistic Regression**
  - **Out: Fri, Oct. 1**
  - **Due: Mon, Oct. 11 at 11:59pm**
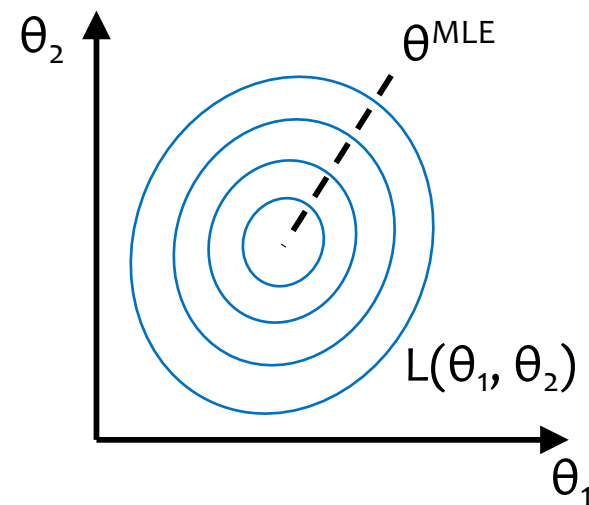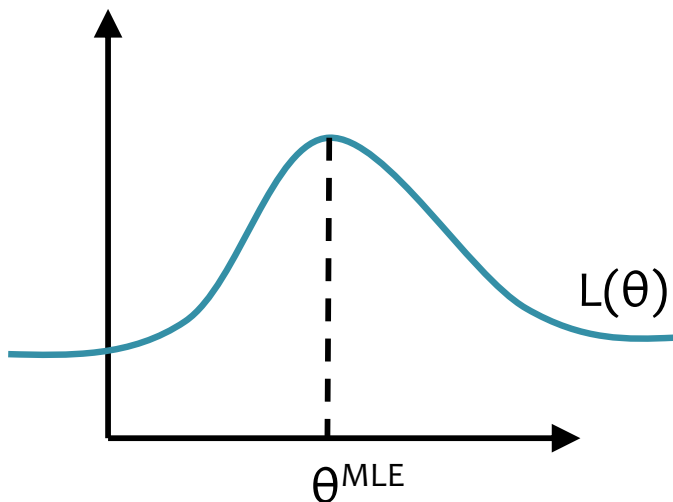
# MAXIMUM LIKELIHOOD ESTIMATION

# MLE

Suppose we have data $\mathcal{D} = \{x^{(i)}\}_{i=1}^N$

**Principle of Maximum Likelihood Estimation:**
Choose the parameters that maximize the likelihood of the data.

$$\boldsymbol{\theta}^{\text{MLE}} = \operatorname*{argmax}_{\boldsymbol{\theta}} \prod_{i=1}^N p(\mathbf{x}^{(i)}|\boldsymbol{\theta})$$

Maximum Likelihood Estimate (MLE)



L(θ)

θ^MLE



θ₂

θ^MLE

L(θ₁, θ₂)

θ₁

# MLE

What does maximizing likelihood accomplish?

- There is only a finite amount of probability mass (i.e. sum-to-one constraint)

- MLE tries to allocate **as much** probability mass **as possible** to the things we have observed…

  …**at the expense** of the things we have **not** observed

# Maximum Likelihood Estimation

The [principle] of Maximum likelihood estimator (MLE):

Choose parameters that make the data "most likely".

<u>Assumptions</u>: Data generated iid from distribution $p^*(x | \vec{\theta}^*)$
and comes from a family of distr parameterized

$$\theta \in \Theta$$

← set of possible parameters

<u>Formally</u>:

$$\theta_{MLE} = \underset{\theta \in \Theta}{\text{argmax}} \; p(D|\theta)$$

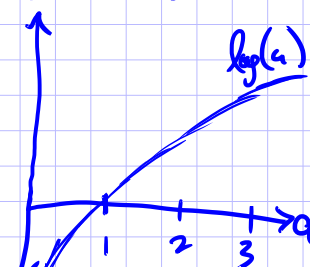$$= \underset{\theta \in \Theta}{\text{argmax}} \; \log p(D|\theta)$$

since log is monotonic

$$= \underset{\theta \in \Theta}{\text{argmax}} \; \ell(\theta)$$

usually a continuos optimization

where $\ell(\theta) \triangleq \log p(D|\theta)$

"log-likelihood"

$\log(a)$

1   2   3   a

$\log(a_1) < \log(a_2)$
iff $a_1 < a_2$
$\Rightarrow \log(f(a_1)) < \log(f(a_2))$
iff $f(a_1) < f(a_2)$

← treat as function of $\theta$
where $D$ is constant

8

# MOTIVATION:
# LOGISTIC REGRESSION

# 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
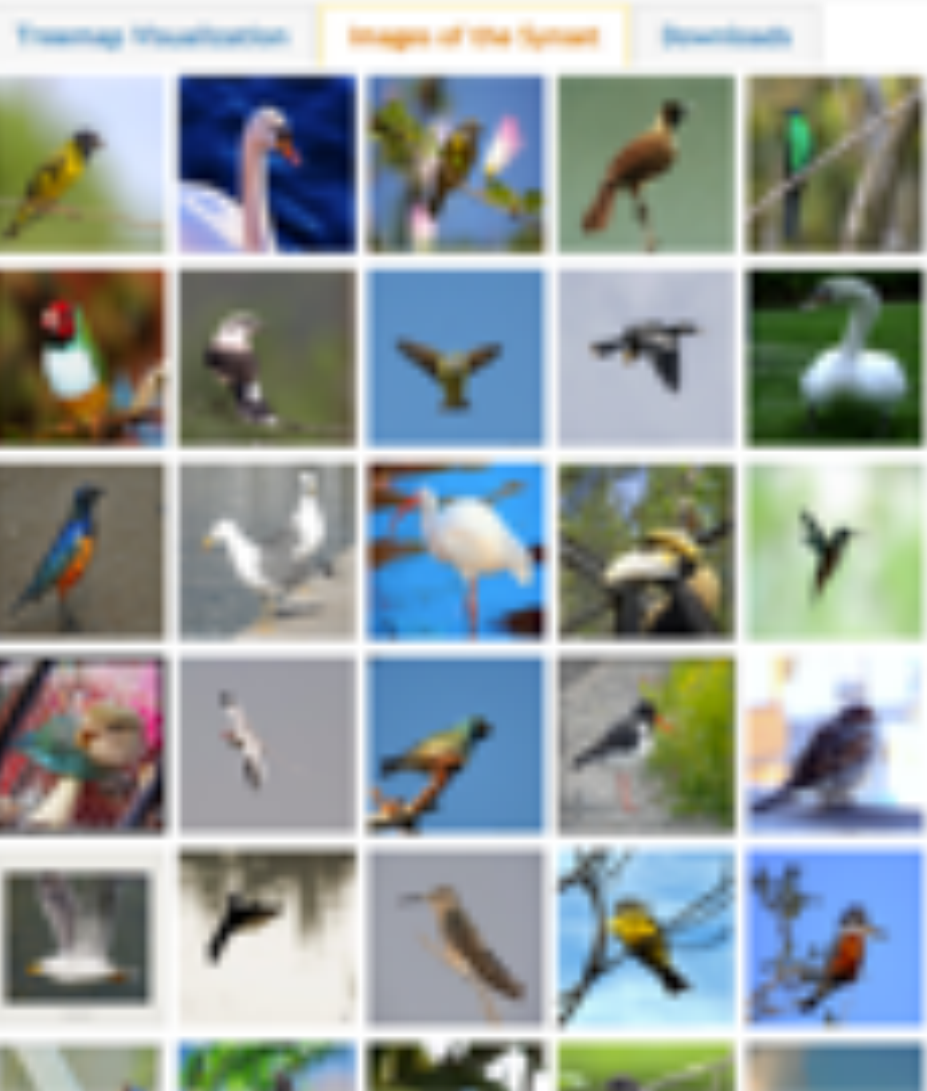  - **Multiclass** classification problem
- Examples from http://image-net.org/

# German iris, Iris kochii

Iris of northern Italy having deep blue purple flowers, similar to but smaller than iris germanica.

# 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

**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

The rest is *just* some fancy feature extraction (discussed later in the course)

This "softmax" layer is Logistic Regression!

224

224

11

11

11

11

Stride of 4

3

48

55

55

Max pooling

128

Max pooling

192

192

3

3

3

3

2048

2048

pooling

2048

2048

dense

dense

1000

# LOGISTIC REGRESSION

# Logistic Regression

**Data:** Inputs are continuous vectors of length M. Outputs are discrete.

$$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N} \text{ where } \mathbf{x} \in \mathbb{R}^M \text{ and } y \in \{0, 1\}$$

We are back to classification.

Despite the name logistic **regression.**

# Linear Models for Classification

Key idea: Try to learn this hyperplane directly

Looking ahead:
- We'll see a number of commonly used Linear Classifiers
- These include:
  – Perceptron
  – Logistic Regression
  – Naïve Bayes (under certain conditions)
  – Support Vector Machines

Directly modeling the hyperplane would use a decision function:

$$h(\mathbf{x}) = \text{sign}(\boldsymbol{\theta}^T \mathbf{x})$$

for:

$$y \in \{-1, +1\}$$

# Background: Hyperplanes

Hyperplane (Definition 1):
$$\mathcal{H} = \{\mathbf{x} : \mathbf{w}^T \mathbf{x} = b\}$$

Hyperplane (Definition 2):
$$\mathcal{H} = \{\mathbf{x'} : \boldsymbol{\theta}^T \mathbf{x'} = 0$$
$$\text{and } x'_1 = 1\}$$
$$\boldsymbol{\theta} = [b, w_1, \ldots, w_M]^T$$
$$\mathbf{x'} = [1, x_1, \ldots, x_M]^T$$

*Notation Trick*: fold the bias *b* and the weights *w* into a single vector **θ** by prepending a constant to *x* and increasing dimensionality by one to get **x'**!

Half-spaces:
$$\mathcal{H}^+ = \{\mathbf{x} : \boldsymbol{\theta}^T \mathbf{x} > 0 \text{ and } x_1 = 1\}$$
$$\mathcal{H}^- = \{\mathbf{x} : \boldsymbol{\theta}^T \mathbf{x} < 0 \text{ and } x_1 = 1\}$$

# Using gradient ascent for linear classifiers

Key idea behind today's lecture:

1. Define a linear classifier (logistic regression)

2. Define an objective function (likelihood)

3. Optimize it with gradient descent to learn parameters

4. Predict the class with highest probability under the model

# Using gradient ascent for linear classifiers

Suppose we wanted to learn a linear classifier, but instead of predicting y ∈ {-1,+1} we wanted to predict y ∈ {0,1}

$$h(\mathbf{x}) = \text{sign}(\boldsymbol{\theta}^T \mathbf{x})$$

sign(u)

# Using gradient ascent for linear classifiers

Suppose we wanted to learn a linear classifier, but instead of predicting y ∈ {-1,+1} we wanted to predict y ∈ {0,1}

$$h(\mathbf{x}) = \text{"sign"}(\boldsymbol{\theta}^T \mathbf{x})$$

"sign"(u)

**Goal:** Learn a linear classifier with Gradient Descent

# Using gradient ascent for linear classifiers

**But this decision function isn't differentiable…**

$$h(\mathbf{x}) = \text{"sign"}(\boldsymbol{\theta}^T \mathbf{x})$$



$$\text{"sign"}(u)$$

**Use a differentiable function instead!**

$$p_{\boldsymbol{\theta}}(y = 1|\mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})}$$



$$\text{logistic}(u) \equiv \frac{1}{1 + e^{-u}}$$

# **Logistic Regression**

*Whiteboard*

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

# Logistic Regression

**Data:** Inputs are continuous vectors of length M. Outputs are discrete.

$$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N} \text{ where } \mathbf{x} \in \mathbb{R}^{M} \text{ and } y \in \{0, 1\}$$

**Model:** Logistic function applied to dot product of parameters with input vector.

$$p_{\boldsymbol{\theta}}(y = 1|\mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^{T}\mathbf{x})}$$

**Learning:** finds the parameters that minimize some objective function.

$$\boldsymbol{\theta}^{*} = \operatorname*{argmin}_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

**Prediction:** Output is the most probable class.

$$\hat{y} = \operatorname*{argmax}_{y \in \{0,1\}} p_{\boldsymbol{\theta}}(y|\mathbf{x})$$

# Logistic Regression

*Whiteboard*

- – Logistic Regression Model
- – Partial derivative for logistic regression
- – Gradient for logistic regression
- – Decision boundary

# LOGISTIC REGRESSION ON GAUSSIAN DATA

# Logistic Regression

# Logistic Regression



Logistic Regression Distribution

# Logistic Regression

## Classification with Logistic Regression

# LEARNING LOGISTIC REGRESSION

# Maximum **Conditional** Likelihood Estimation

**Learning:** finds the parameters that minimize some objective function**.**

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \ J(\boldsymbol{\theta})$$

We minimize the *negative* log conditional likelihood:

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

Why?

1. We can't maximize likelihood (as in Naïve Bayes) because we don't have a joint model p(x,y)
2. It worked well for Linear Regression (least squares is MCLE)

# Maximum **Conditional** Likelihood Estimation

**Learning:** Four approaches to solving $\theta^* = \underset{\theta}{\mathrm{argmin}}\, J(\theta)$

**Approach 1:** Gradient Descent
(take larger – more certain – steps opposite the gradient)

**Approach 2:** Stochastic Gradient Descent (SGD)
(take many small steps opposite the gradient)

**Approach 3:** Newton's Method
(use second derivatives to better follow curvature)

**Approach 4:** Closed Form???
(set derivatives equal to zero and solve for parameters)

# Maximum **Conditional** Likelihood Estimation

**Learning:** Four approaches to solving $\theta^* = \underset{\theta}{\operatorname{argmin}} J(\theta)$

**Approach 1:** Gradient Descent
(take larger – more certain – steps opposite the gradient)

**Approach 2:** Stochastic Gradient Descent (SGD)
(take many small steps opposite the gradient)

**Approach 3:** Newton's Method
(use second derivatives to better follow curvature)

~~**Approach 4:** Closed Form???~~
(set derivatives equal to zero and solve for parameters)

Logistic Regression does not have a closed form solution for MLE parameters.

# 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

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

E. (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

# Lecture 10: In-Class Poll



**0 done**

# Question 1

| A |
|---|
| B |
| C |
| D |
| E |
| F |

# Gradient Descent

**Algorithm 1** Gradient Descent

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

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** $\boldsymbol{\theta}$



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

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \begin{bmatrix} \frac{d}{d\theta_1} J(\boldsymbol{\theta}) \\ \frac{d}{d\theta_2} J(\boldsymbol{\theta}) \\ \vdots \\ \frac{d}{d\theta_M} J(\boldsymbol{\theta}) \end{bmatrix}$$

# 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** $i \in$ shuffle($\{1, 2, \ldots, N\}$) **do**
5:             $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \gamma \nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta})$
6:     **return** $\theta$

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

We need a per-example objective:

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

43

# Logistic Regression vs. Perceptron

**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:**

# Question 2

A

B

C

# 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: $\theta \leftarrow \theta - \lambda g$

## Three variants of first-order optimization:

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

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

Mini-batch SGD: $g = \dfrac{1}{S} \displaystyle\sum_{s=1}^{S} \nabla J^{(i_s)}(\boldsymbol{\theta})$      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

$$p(y|x) \propto \exp(\Theta_y \bullet f( \quad ))$$

# Where do features come from?



Feature Engineering

hand-crafted features

Sun et al., 2011

Zhou et al., 2005

*First word before M1*
*Second word before M1*
*Bag-of-words in M1*
*Head word of M1*
*Other word in between*
*First word after M2*
*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*
*WordNet Tags*
*Heads of chunks in between*
*Path of phrase labels*
*Combination of entity types*

Feature Learning

# Where do features come from?

# Where do features come from?



Feature Engineering

pooling

The [movie] showed [wars]

**Convolutional Neural Networks**
*(Collobert and Weston 2008)*

**CNN**

The [movie] showed [wars]

**Recursive Auto Encoder**
*(Socher 2011)*

**RAE**

Zhou et al., 2005

**word embeddings**
Mikolov et al., 2013

**string embeddings**
Socher, 2011
Collobert & Weston, 2008

Feature Learning

# Where do features come from?

# Where do features come from?



Feature Engineering

Feature Learning

hand-crafted features

Sun et al., 2011

Zhou et al., 2005

word embedding features

Koo et al. 2008

Turian et al. 2010

Hermann et al. 2014

Refine embedding features with semantic/syntactic info

tree embeddings
Socher et al., 2013
Hermann & Blunsom, 2013

word embeddings
Mikolov et al., 2013

string embeddings
Socher, 2011
Collobert & Weston, 2008

66

# Where do features come from?

# Feature Engineering for NLP

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?**

| deter. | noun | noun | verb | verb | noun |
|--------|------|------|------|------|------|
| *The* | *movie* | *I* | *watched* | *depicted* | *hope* |

# Feature Engineering for NLP

**Per-word Features:**

|  | $x^{(1)}$ | $x^{(2)}$ | $x^{(3)}$ | $x^{(4)}$ | $x^{(5)}$ | $x^{(6)}$ |
|---|---|---|---|---|---|---|
| `is-capital(wi)` | 1 | 0 | 1 | 0 | 0 | 0 |
| `endswith(wi,"e")` | 1 | 1 | 0 | 0 | 0 | 1 |
| `endswith(wi,"d")` | 0 | 0 | 0 | 1 | 1 | 0 |
| `endswith(wi,"ed")` | 0 | 0 | 0 | 1 | 1 | 0 |
| `wi == "aardvark"` | 0 | 0 | 0 | 0 | 0 | 0 |
| `wi == "hope"` | 0 | 0 | 0 | 0 | 0 | 1 |
| ... | ... | ... | ... | ... | ... | ... |

| deter. | noun | noun | verb | verb | noun |
|---|---|---|---|---|---|
| *The* | *movie* | *I* | *watched* | *depicted* | *hope* |

69

# Feature Engineering for NLP

**Context Features:**

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

| deter. | noun | noun | verb | verb | noun |
|---|---|---|---|---|---|
| *The* | *movie* | *I* | *watched* | *depicted* | *hope* |

# Feature Engineering for NLP

**Context Features:**

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

| deter. | noun | noun | verb | verb | noun |
|---|---|---|---|---|---|
| *The* | *movie* | *I* | *watched* | *depicted* | *hope* |

# Feature Engineering for NLP

**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* |

# Feature Engineering for CV

Edge detection (Canny)



Corner Detection (Harris)

Figures from http://opencv.org

# Feature Engineering for CV

## Scale Invariant Feature Transform (SIFT)



Figure 3: Model images of planar objects are shown in the top row. Recognition results below show model outlines and image 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.

Figure from Lowe (1999) and Lowe (2004)

# NON-LINEAR FEATURES

# Nonlinear Features

- aka. "nonlinear basis functions"
- So far, input was always $\mathbf{x} = [x_1, \ldots, x_M]$
- **Key Idea**: let input be some function of **x**
  - original input: $\mathbf{x} \in \mathbb{R}^M$     where $M' > M$ (usually)
  - new input: $\mathbf{x}' \in \mathbb{R}^{M'}$
  - define $\mathbf{x}' = b(\mathbf{x}) = [b_1(\mathbf{x}), b_2(\mathbf{x}), \ldots, b_{M'}(\mathbf{x})]$

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

- **Examples:** ($M = 1$)

| | |
|---|---|
| polynomial | $b_j(x) = x^j \quad \forall j \in \{1, \ldots, J\}$ |
| radial basis function | $b_j(x) = \exp\left(\dfrac{-(x - \mu_j)^2}{2\sigma_j^2}\right)$ |
| sigmoid | $b_j(x) = \dfrac{1}{1 + \exp(-\omega_j x)}$ |
| 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 $y = \mathbf{w}^T f(\mathbf{x}) + b$ where f(.) is a polynomial basis function

| i | y | x |
|---|---|---|
| 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 | y | x |
|---|---|---|
| 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



Linear Regression (poly=1)

# Example: Linear Regression

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

| i | y | x | x² |
|---|---|---|---|
| 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



Linear Regression (poly=2)

# Example: Linear Regression

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

| i | y | x | x² | x³ |
|---|---|---|---|---|
| 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 | $(1.9)^2$ | $(1.9)^3$ |

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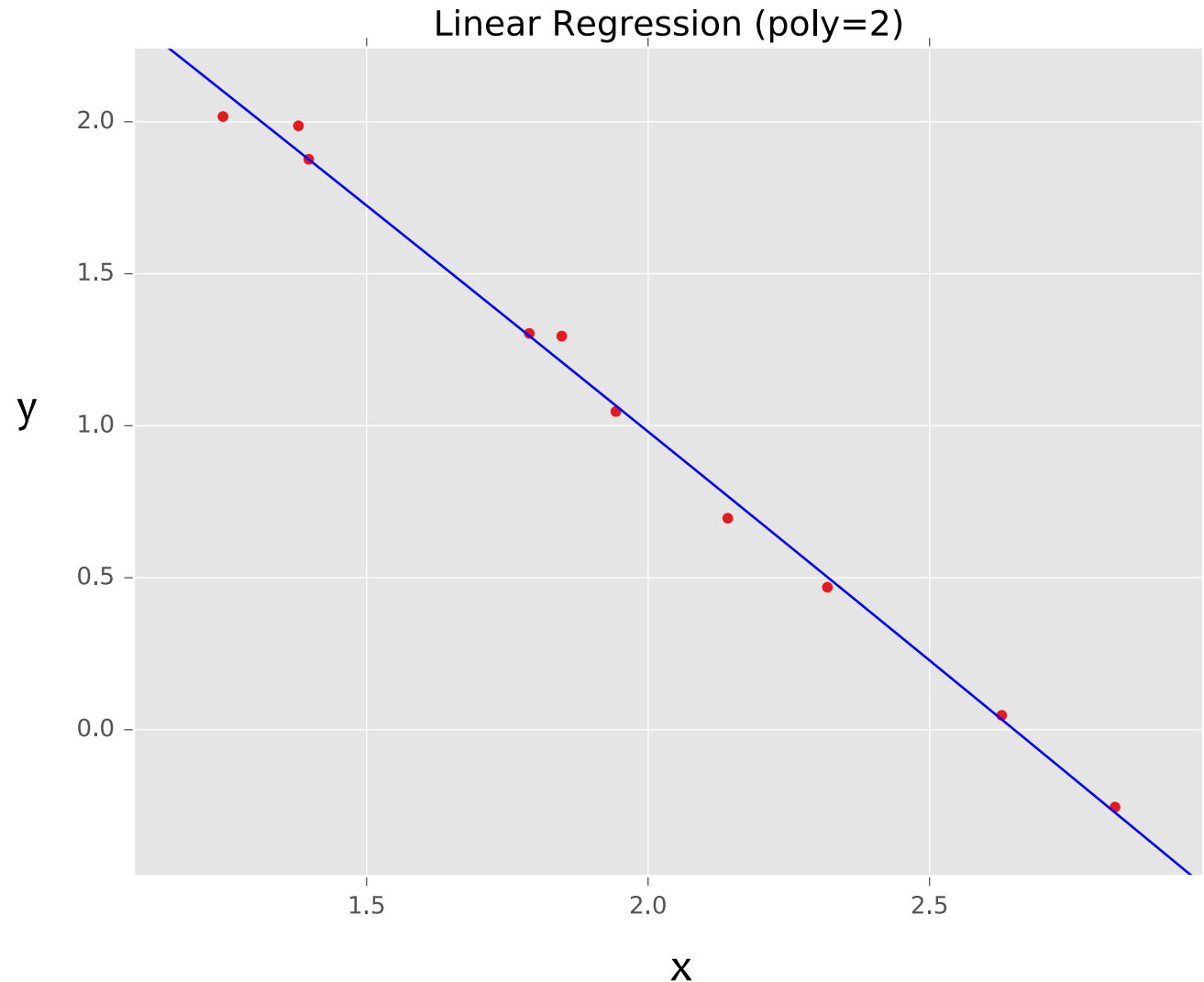
y



Linear Regression (poly=3)

x

# Example: Linear Regression

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

| i | y | x | ... | $x^5$ |
|---|---|---|---|---|
| 1 | 2.0 | 1.2 | ... | $(1.2)^5$ |
| 2 | 1.3 | 1.7 | ... | $(1.7)^5$ |
| ... | ... | ... | ... | ... |
| 10 | 1.1 | 1.9 | ... | $(1.9)^5$ |

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



Linear Regression (poly=5)

y

x

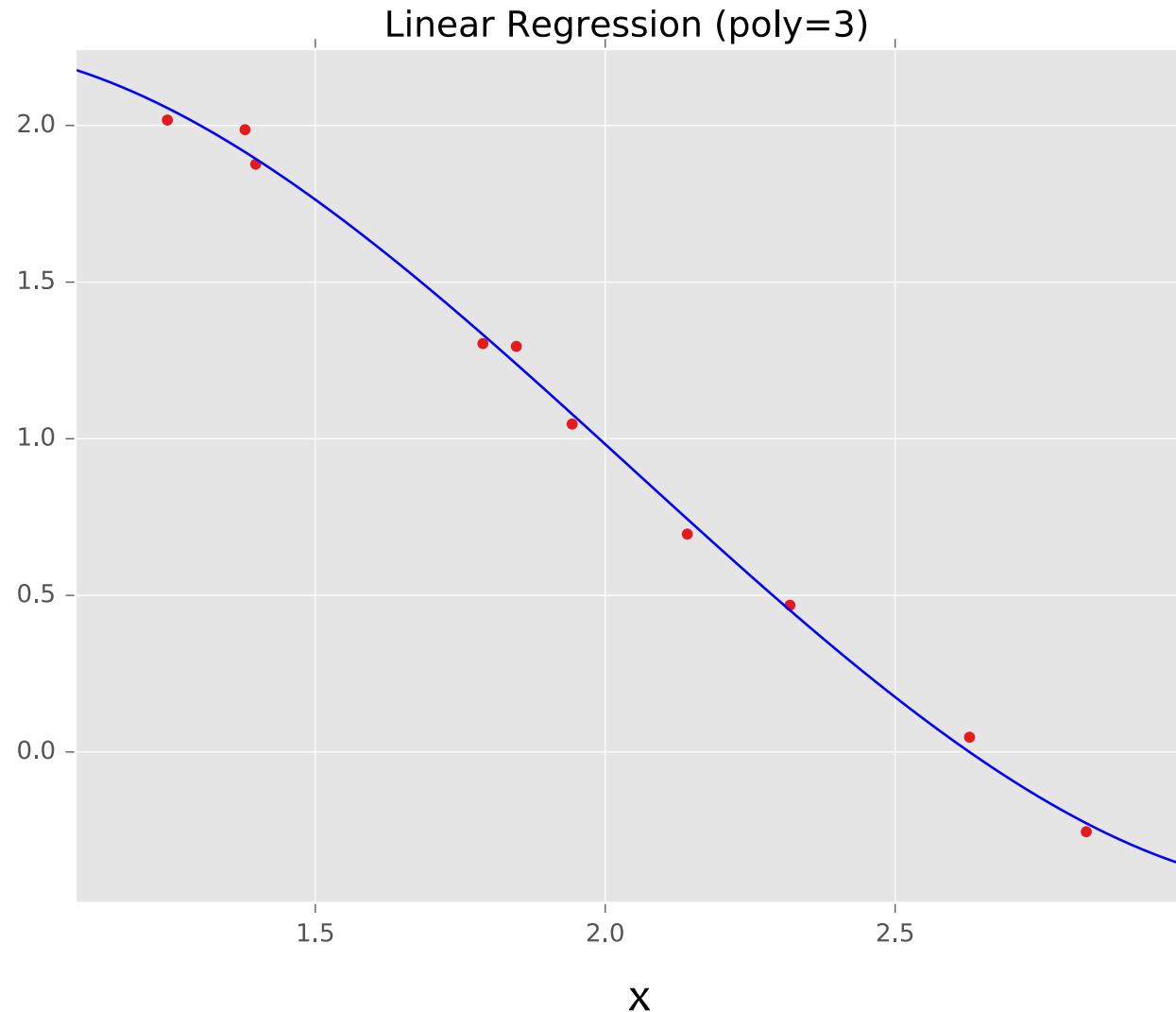# Example: Linear Regression

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

| i | y | x | ... | $x^8$ |
|---|---|---|---|---|
| 1 | 2.0 | 1.2 | ... | $(1.2)^8$ |
| 2 | 1.3 | 1.7 | ... | $(1.7)^8$ |
| ... | ... | ... | ... | ... |
| 10 | 1.1 | 1.9 | ... | $(1.9)^8$ |

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



Linear Regression (poly=8)

# Example: Linear Regression

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

| i | y | x | ... | $x^9$ |
|---|---|---|---|---|
| 1 | 2.0 | 1.2 | ... | $(1.2)^9$ |
| 2 | 1.3 | 1.7 | ... | $(1.7)^9$ |
| ... | ... | ... | ... | ... |
| 10 | 1.1 | 1.9 | ... | $(1.9)^9$ |

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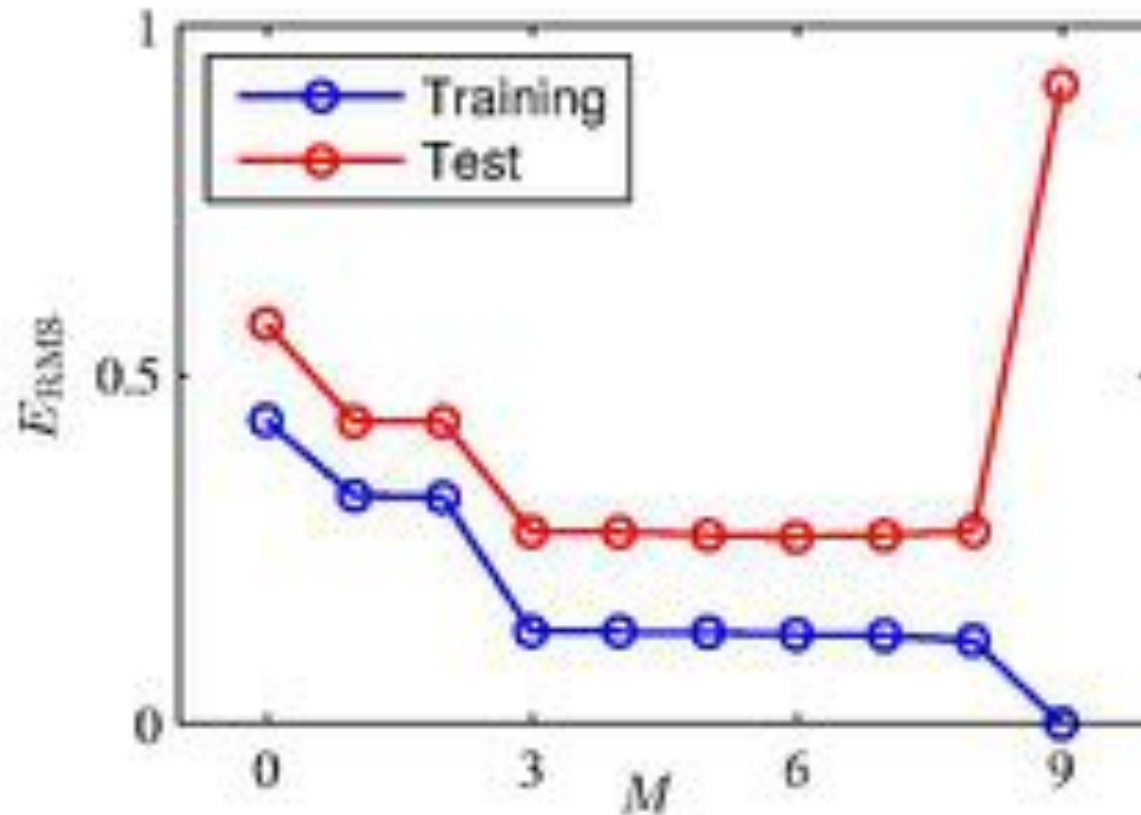


Linear Regression (poly=9)

# Over-fitting



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

# Polynomial Coefficients

|            | $M = 0$ | $M = 1$ | $M = 3$ | $M = 9$      |
|------------|---------|---------|---------|--------------|
| $\theta_0$ | 0.19    | 0.82    | 0.31    | 0.35         |
| $\theta_1$ |         | -1.27   | 7.99    | 232.37       |
| $\theta_2$ |         |         | -25.43  | -5321.83     |
| $\theta_3$ |         |         | 17.37   | 48568.31     |
| $\theta_4$ |         |         |         | -231639.30   |
| $\theta_5$ |         |         |         | 640042.26    |
| $\theta_6$ |         |         |         | -1061800.52  |
| $\theta_7$ |         |         |         | 1042400.18   |
| $\theta_8$ |         |         |         | -557682.99   |
| $\theta_9$ |         |         |         | 125201.43    |

# Example: Linear Regression

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

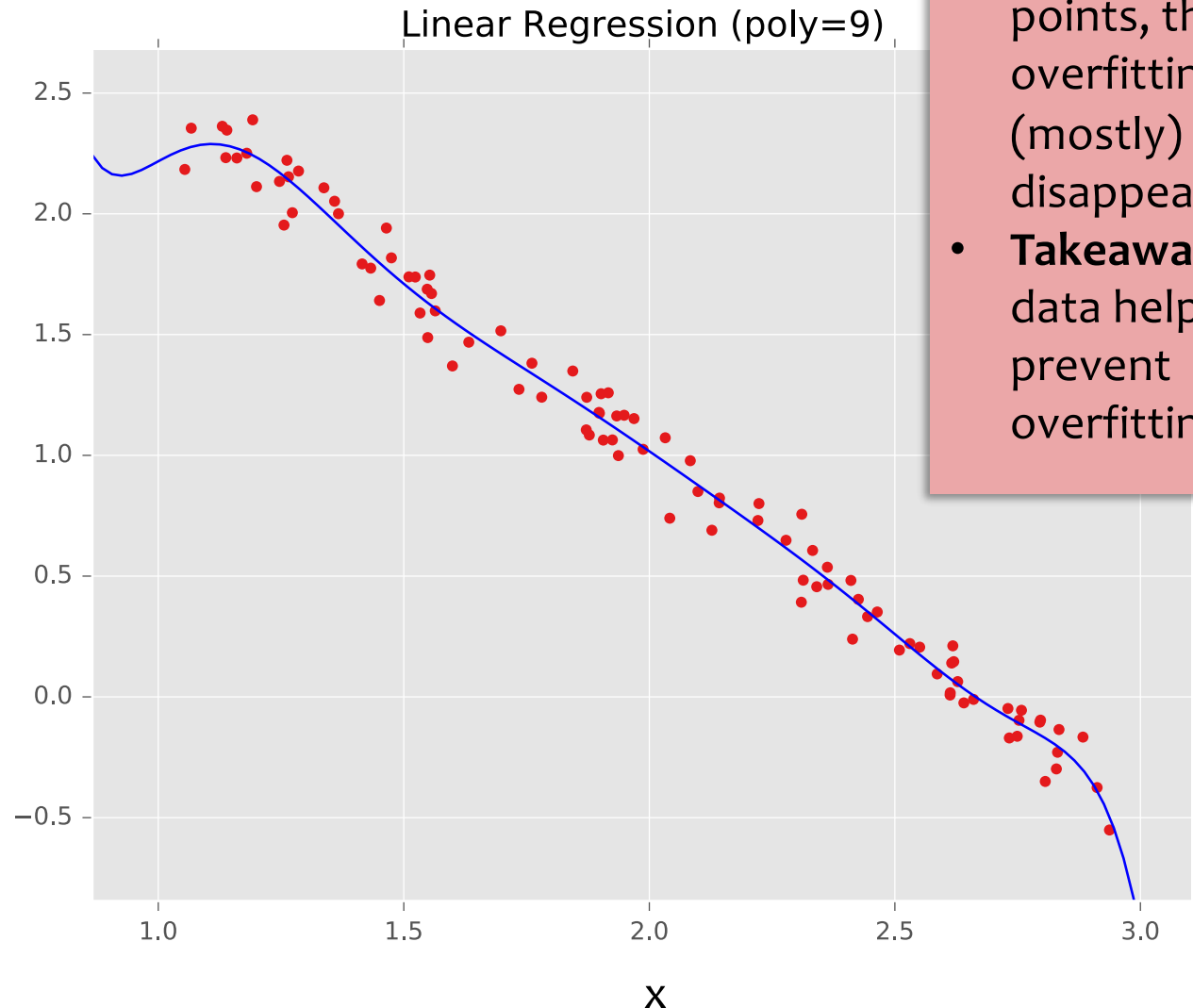| i | y | x | ... | $x^9$ |
|---|---|---|-----|-------|
| 1 | 2.0 | 1.2 | ... | $(1.2)^9$ |
| 2 | 1.3 | 1.7 | ... | $(1.7)^9$ |
| ... | ... | ... | ... | ... |
| 10 | 1.1 | 1.9 | ... | $(1.9)^9$ |

y

Linear Regression (poly=9)

x

- With just N = 10 points we overfit!
- But with N = 100 points, the overfitting (mostly) disappears
- **Takeaway:** more data helps prevent overfitting

# Example: Linear Regression

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

| i | y | x | ... | $x^9$ |
|---|---|---|-----|-------|
| 1 | 2.0 | 1.2 | ... | $(1.2)^9$ |
| 2 | 1.3 | 1.7 | ... | $(1.7)^9$ |
| 3 | 0.1 | 2.7 | ... | $(2.7)^9$ |
| 4 | 1.1 | 1.9 | ... | $(1.9)^9$ |
| ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |
| 98 | ... | ... | ... | ... |
| 99 | ... | ... | ... | ... |
| 100 | 0.9 | 1.5 | ... | $(1.5)^9$ |



Linear Regression (poly=9)

- With just N = 10 points we overfit!
- But with N = 100 points, the overfitting (mostly) disappears
- **Takeaway:** more data helps prevent overfitting