

Logistic regression

- A general classifier (features can be anything, not just word counts)
- Discriminative instead of generative
- Building block for neural networks

Generative vs. discriminative classifiers

- The goal of any classifier is to compute $P(c | d)$ for all $c \in C$
- Or just $\arg \max_{c \in C} P(c | d)$
- A generative model gives us $P(d | c)$
- Generative classifier:
$$P(c | d) = \frac{P(d | c)P(c)}{\sum_{c' \in C} P(d | c')P(c')}$$

Discriminative classifier

- Discriminative classifiers model $P(c | d)$ directly
- For each instance (example) d , we have features $\mathbf{x} = (x_1, \dots, x_n)$
- This is all we know about the document: we are going to compute $P(c | \mathbf{x})$
- Easier: we don't need to know how to compute $P(x_n | x_1, \dots, x_{n-1})$

Building blocks of a discriminative classifier

- Features: $d \rightarrow \mathbf{x}$
- Classification function $P(c | \mathbf{x}) = \sigma(\mathbf{w} \cdot \mathbf{x} + b)$
- Loss function: cross-entropy
- Optimization method: stochastic gradient descent

Classification function

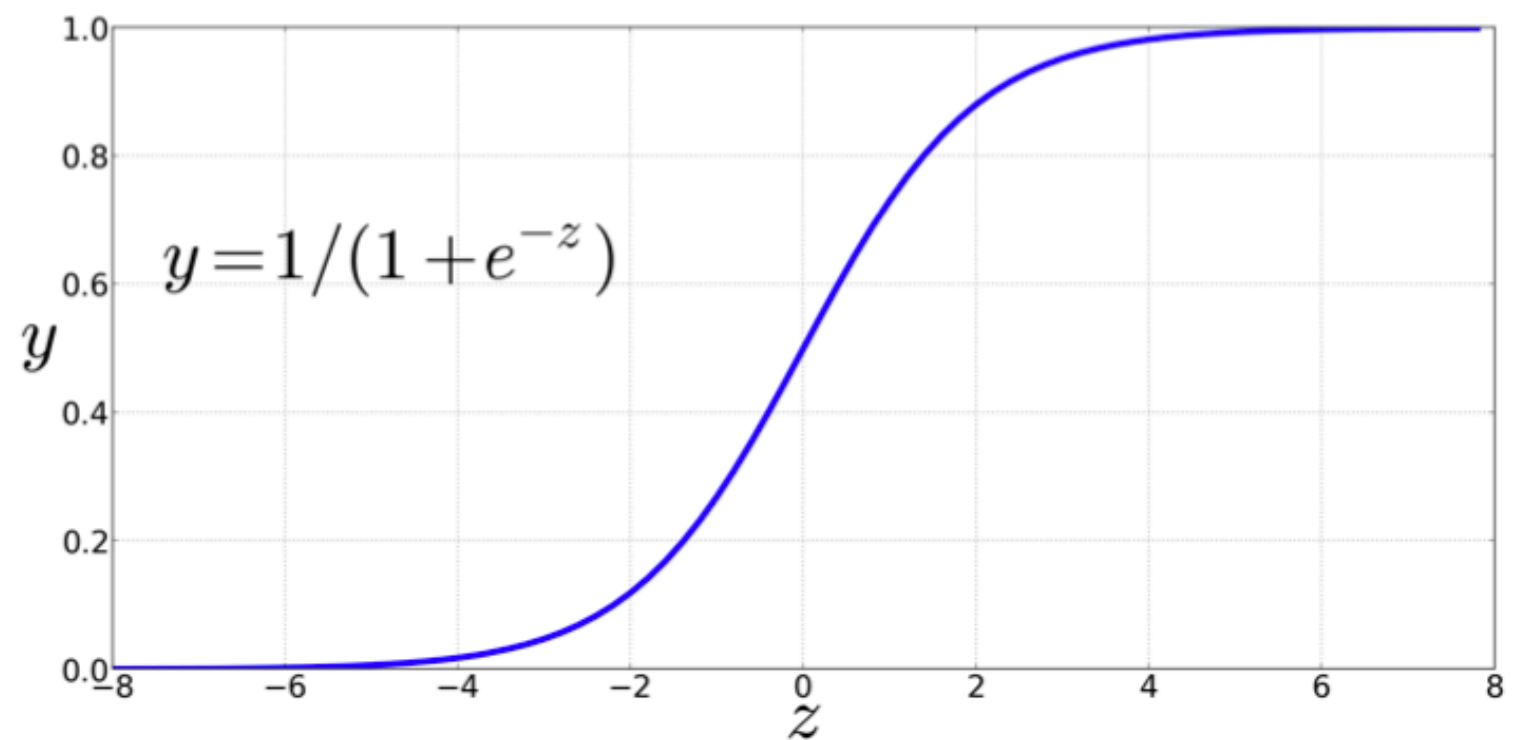
$$z = \mathbf{w} \cdot \mathbf{x} + b = \sum_{j=1}^n w_j x_j + b$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

$$P(y = 1) = \hat{y}$$

$$P(y = 0) = 1 - \hat{y}$$

$$P(y = 1) > 0.5?$$



Example

Var	Definition	Value in Fig. 5.2
x_1	count(positive lexicon) \in doc)	3
x_2	count(negative lexicon) \in doc)	2
x_3	$\begin{cases} 1 & \text{if "no" } \in \text{ doc} \\ 0 & \text{otherwise} \end{cases}$	1
x_4	count(1st and 2nd pronouns \in doc)	3
x_5	$\begin{cases} 1 & \text{if "!" } \in \text{ doc} \\ 0 & \text{otherwise} \end{cases}$	0
x_6	log(word count of doc)	$\ln(66) = 4.19$

It's **hokey**. There are virtually **no** surprises, and the writing is **second-rate**. So why was it so **enjoyable**? For one thing, the cast is **great**. Another **nice** touch is the music. **I** was overcome with the urge to get off the couch and start dancing. It sucked **me** in, and it'll do the same to **you**.

$x_2=2$ $x_3=1$ $x_1=3$ $x_5=0$ $x_6=4.19$ $x_4=3$

Loss function

$$\mathbf{x}^{(i)}, y^{(i)}, \hat{y}^{(i)}$$

$$L_{\theta}(\hat{y}^{(i)}, y^{(i)}) \quad \theta = (\mathbf{w}, b)$$

$$\arg \min_{\theta} \sum_{i=1}^k L_{\theta}(\hat{y}^{(i)}, y^{(i)})$$

Logistic loss function

$$\hat{y} = \sigma(\mathbf{w} \cdot \mathbf{x} + b)$$

$$\hat{P}(y = 1 | x) = \hat{y}$$

$$\hat{P}(y = 0 | x) = 1 - \hat{y}$$

$$\hat{P}(y | x) = \hat{y}^y (1 - \hat{y})^{1-y}$$

$$\log \hat{P}(y | x) = y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})$$

$$L(\hat{y}, y) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

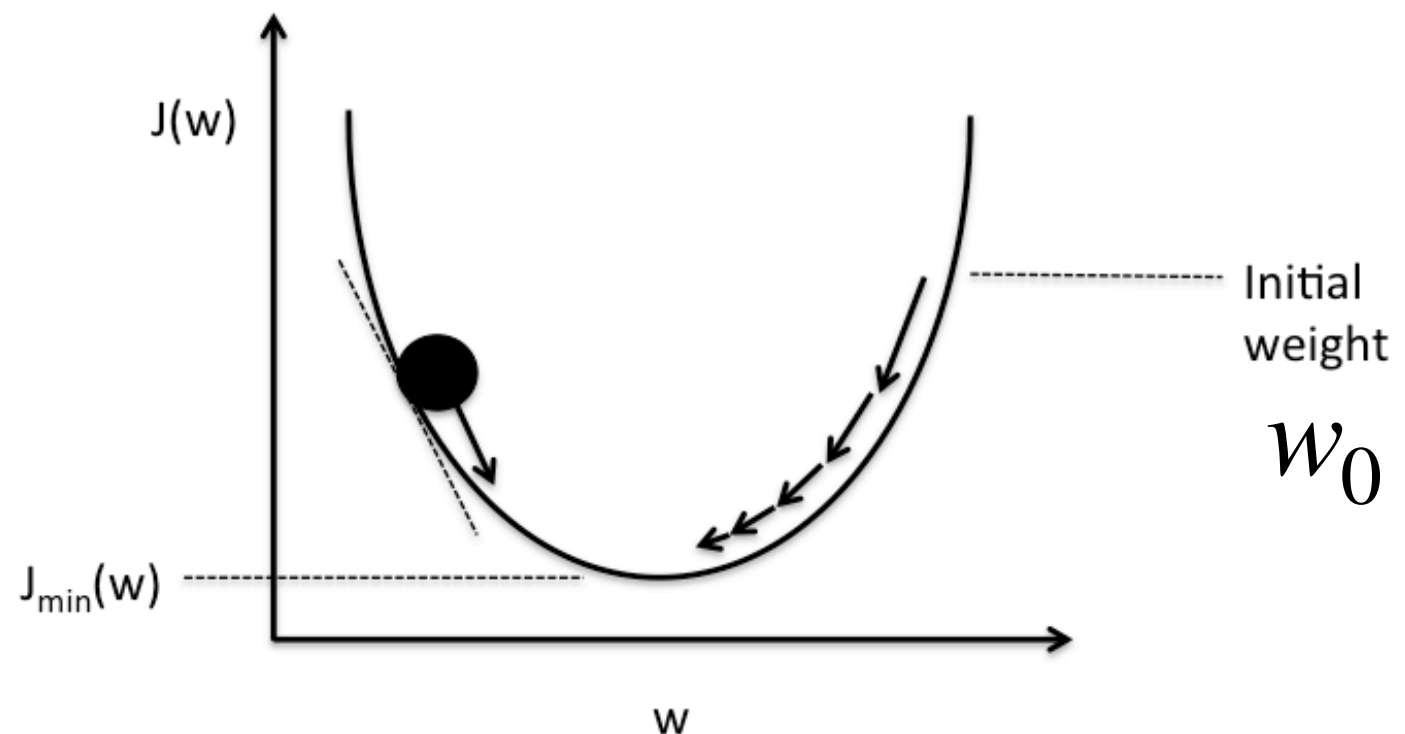
Gradient descent

- Goal: find the weights that minimize the loss
 $L(\hat{y}, y) = L(w)$
- The derivative of the loss indicates whether the loss is increasing or decreasing in w , and at what rate

The derivative of the
loss function

$$w^{t+1} = w^t - \eta L'(w)$$

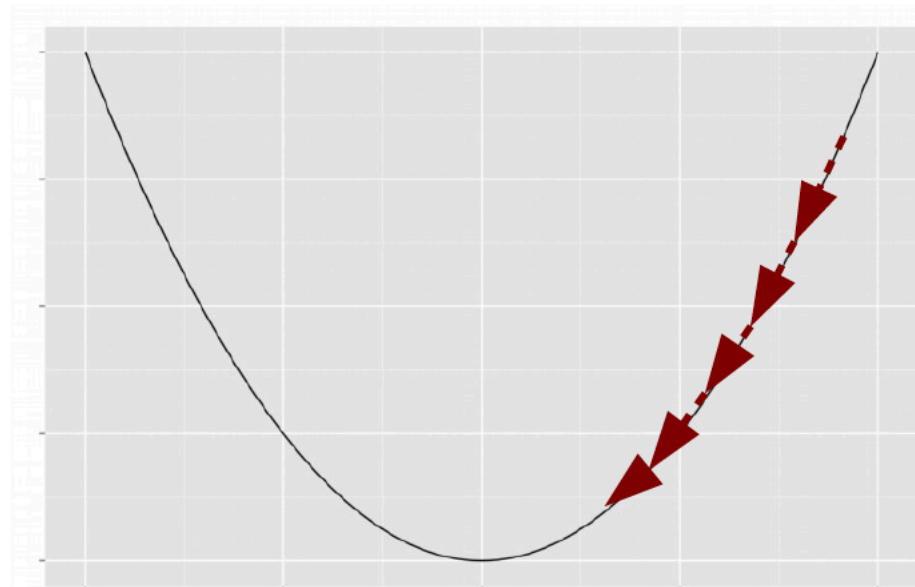
Learning rate



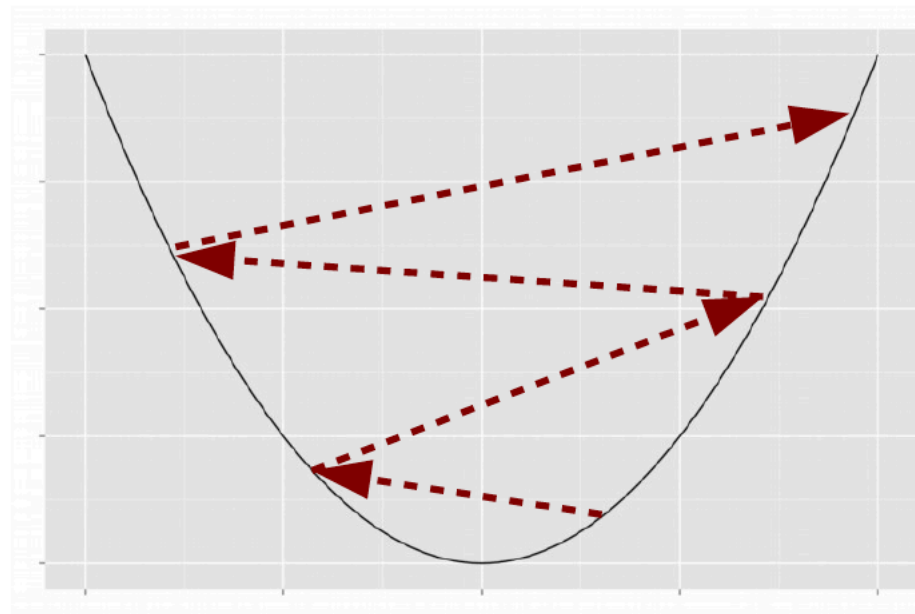
Schematic of gradient descent.

Gradient descent: learning rate

Lower learning rate:
slow convergence



Higher learning
rate: may not
always converge

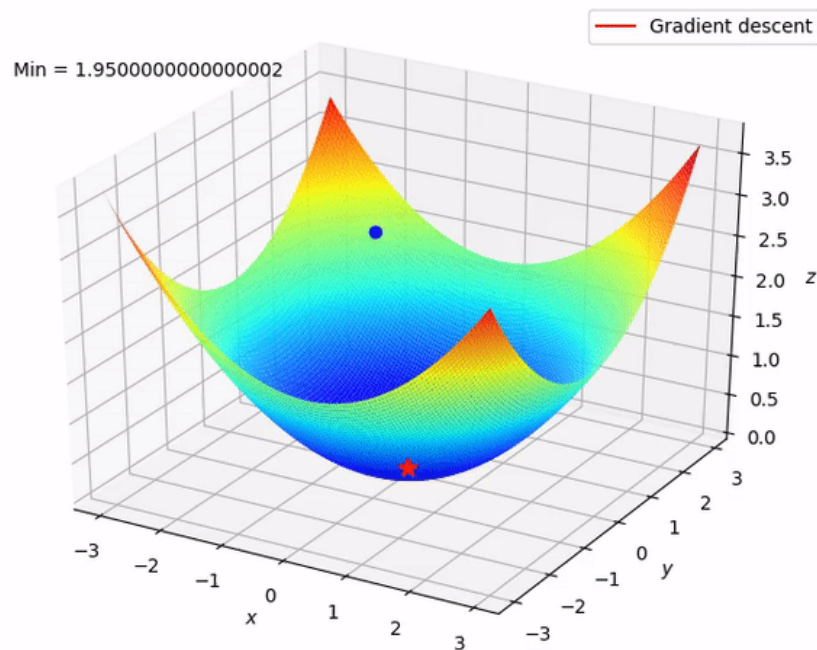


Gradient descent

- With multiple variables: move in the direction of the gradient of the loss function

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla L(\mathbf{w})$$

$$\nabla L(\mathbf{w}) = \left(\frac{\partial L}{\partial w_1}, \dots, \frac{\partial J}{\partial w_n} \right)$$



Gradient descent for logistic regression

Loss for the instance $(\mathbf{x}^{(i)}, y^{(i)})$:

$$L^{(i)}(\mathbf{w}, b) = -y^{(i)} \log(\sigma(\mathbf{w} \cdot \mathbf{x}^{(i)} + b)) - (1 - y^{(i)}) \log(1 - \sigma(\mathbf{w} \cdot \mathbf{x}^{(i)} + b))$$

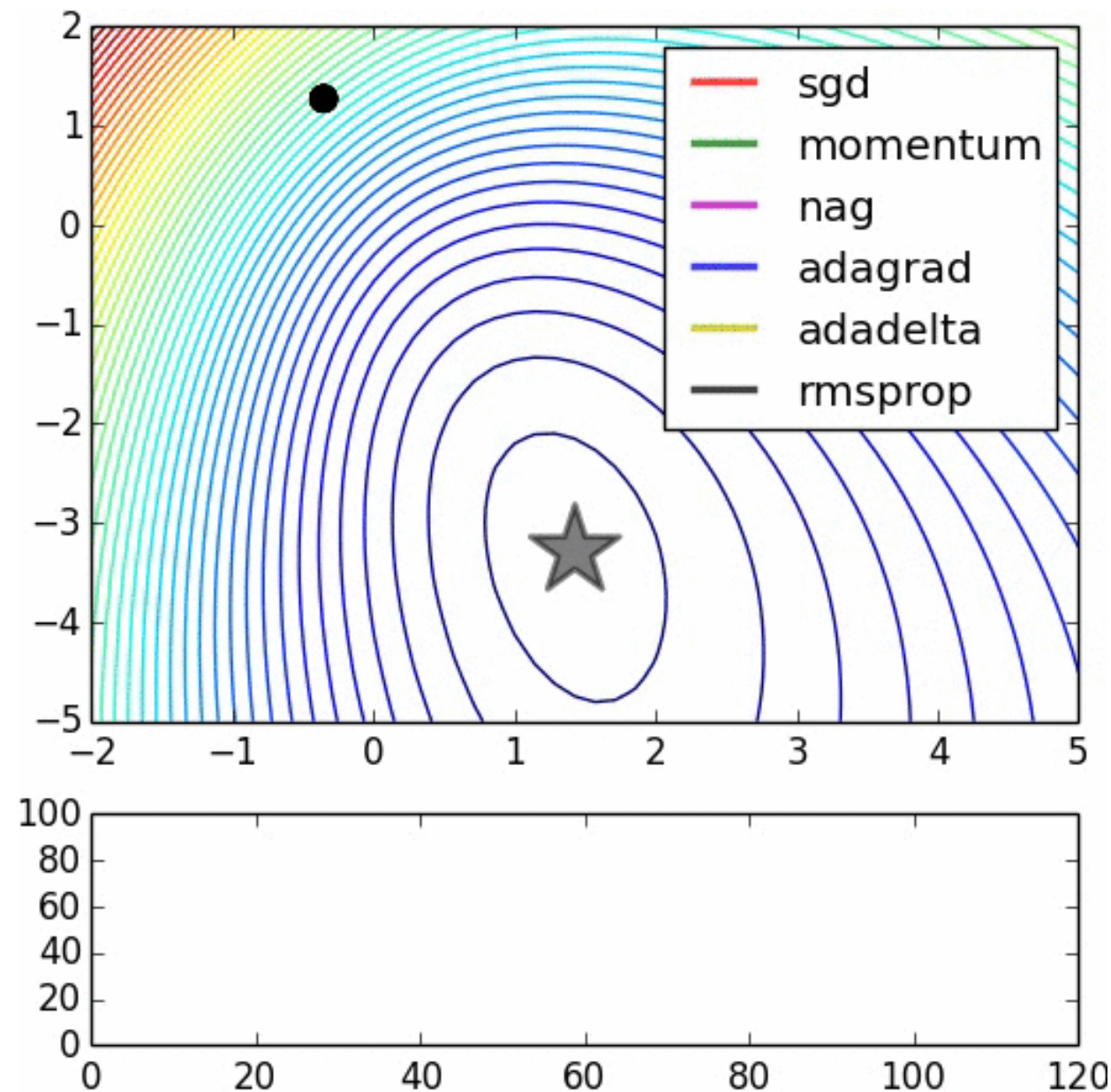
$$\frac{\partial L^{(i)}(\mathbf{w}, b)}{\partial w_j} = (\sigma(\mathbf{w} \cdot \mathbf{x}^{(i)} + b) - y^{(i)}) x_j^{(i)}$$

$$L(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m L^{(i)}(\mathbf{w}, b)$$

Stochastic gradient descent

- The loss is a function of all the training instances
- Stochastic gradient descent (SGD): compute the loss gradient after each example, and move in that direction
- Noisy estimate of the gradient over the full training data: mini-batching can reduce the variance of this estimate

Variations on stochastic gradient descent



(Credit: Alec Radford)

Skip-gram with negative sampling (SGNS)

... lemon, a [tablespoon of apricot jam, a] pinch ...
 c1 c2 t c3 c4

$$P(+ | t, c) = \sigma(\mathbf{t} \cdot \mathbf{c}) \qquad P(- | t, c) = 1 - \sigma(\mathbf{t} \cdot \mathbf{c})$$

$$P(+ | t, c_1, \dots, c_k) = \prod_{i=1}^k P(+ | t, c_i)$$

Skip-gram with negative sampling (SGNS)

... lemon, a [tablespoon of apricot jam, a] pinch ...
 c1 c2 t c3 c4

$$P(+ | t, c) = \sigma(\mathbf{t} \cdot \mathbf{c})$$

$$P(- | t, c) = 1 - \sigma(\mathbf{t} \cdot \mathbf{c})$$

positive examples +

t	c
---	---

apricot	tablespoon
apricot	of
apricot	jam
apricot	a

negative examples -

t	c	t	c
---	---	---	---

apricot	aardvark	apricot	seven
apricot	my	apricot	forever
apricot	where	apricot	dear
apricot	coaxial	apricot	if

Skip-gram with negative sampling (SGNS)

$$\theta = (\mathbf{t}_{\text{dog}}, \mathbf{c}_{\text{dog}}, \mathbf{t}_{\text{penguin}}, \mathbf{c}_{\text{penguin}}, \dots)$$

$$F(\theta) = \sum_{(t,c) \in +} \log(P(+ | t, c)) + \sum_{(t,c) \in -} \log(P(- | t, c))$$

$$P(+ | t, c) = \sigma(\mathbf{t} \cdot \mathbf{c}) \quad P(- | t, c) = 1 - \sigma(\mathbf{t} \cdot \mathbf{c})$$