

Log-Linear Models

Kevin Duh
Intro to NLP, Fall 2019

Motivation

- We want to **estimate $p(y|x)$**
 - e.g. $P(Y=\text{word} \mid X = \text{previous_words})$
 - e.g. $P(Y=\text{label} \mid X = \text{document})$
- In the n-gram lectures, we estimated this by counts:
 - $P(Y=\text{word} \mid X = \text{previous words}) = \frac{\text{Count}(\text{word}, \text{previous_words})}{\text{Count}(\text{previous_words})}$
 - Smoothing alleviates unreliable estimates

Outline

1. What is a Log-Linear Model and why use it?
2. How to train it?
3. Interactive visualization

Motivation

- A different approach to estimate $p(y|x)$
- Log-linear model:
 - enables us to incorporate our knowledge of the problem by defining features of x,y

Log-Linear Model Definition

$$p(y|x) \propto \text{score}(x, y)$$

$$\text{score}(x, y) = \sum_k \theta_k f_k(x, y)$$

1. Define K **feature functions** $f()$, each measuring something about x and y
2. The probability depends on a **weighted combination** of $f()$
3. **Weights are learned** from training data

Example features for word prediction $p(w_n | w_{n-1}, w_{n-2})$

$$P(W_n = \text{"Sam"} | W_{n-1} = \text{"am"}, W_{n-2} = \text{"I"})$$

$f_1(x, y) =$	$Count(\text{"I am Sam"})$
$f_2(x, y) =$	$1/Count(\text{"I am"})$
$f_3(x, y) =$	$Count(\text{"am Sam"})$
$f_4(x, y) =$	$1/Count(\text{"am"})$
$f_5(x, y) =$	$Count(\text{"Sam"})$
$f_6(x, y) =$	$P_{add-one-smooth}(\text{"Sam"} \text{"I am"})$
$f_7(x, y) =$	$P_{add-two-smooth}(\text{"Sam"} \text{"I am"})$
$f_8(x, y) =$	$P_{add-one-smooth}(\text{"Sam"} \text{"I"})$
$f_9(x, y) =$	$I(\text{"Sam"} \text{ is capitalized})$
$f_{10}(x, y) =$	$I(\text{"Sam"} \text{ is Noun and previous word is Verb})$

Log-Linear Model Definition

We said $p(y|x) \propto \text{score}(x, y)$

Main Concept!
Get familiar with it

More precisely, Log-Linear models are defined as:

$$p(y | x) = \frac{1}{Z(x)} \exp(\text{score}(x, y)) = \frac{1}{Z(x)} \exp\left(\sum_{k=1}^K \theta_k f(x, y)\right)$$

Normalization:

$$Z(x) = \sum_{y'} \exp(\text{score}(x, y'))$$

It's called Log-Linear because $\log(p(y|x))$ is a linear function.

Sometimes also called Maximum Entropy (MaxEnt) or Multinomial Logistic Regression

Outline

1. What is a Log-Linear Model and why use it?
2. How to train it?
3. Interactive visualization

Training the parameters

We're given training data: $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

Want to find parameters that maximize likelihood on training data:

$$p(\text{training data}; \vec{\theta}) = \prod_{n=1}^N p(y_n | x_n; \vec{\theta})$$

$$p(y_n | x_n; \vec{\theta}) = \frac{1}{Z(x)} \exp\left(\sum_{k=1}^K \theta_k f(x_n, y_n)\right)$$
$$\vec{\theta} = [\theta_1, \theta_2, \dots, \theta_K]$$

Equivalent to maximizing the log-likelihood:

$$\sum_{n=1}^N \log(p(y_n | x_n; \vec{\theta}))$$

Side-note: Probability vs. Likelihood

- **Parameter known. Data unknown.**
 - I know a coin is biased with this parameter: $P(\text{Flip}=\text{H})=0.7$, $P(\text{Flip}=\text{T})=0.3$
 - Question: If I flip 3 times, what is the probability I'll get HHH? $0.7 \times 0.7 \times 0.7 = 0.343$
- **Data known. Parameter unknown.**
 - I know that I flipped 3 times and got HHH.
 - Question: What's my estimate of the biasness of the coin? i.e. $P(\text{Flip}=\text{H})=???$
 - Maximum likelihood solution is $P(\text{Flip}=\text{H})=1.0$, $P(\text{Flip}=\text{T})=0.0$.
 - Try $P(\text{Flip}=\text{H})=1$. Likelihood = $1.0 \times 1.0 \times 1.0 = 1.0$
 - Try $P(\text{Flip}=\text{H})=0.7$. Likelihood = $0.7 \times 0.7 \times 0.7 = 0.343$

Regularization Term

- Large weights may lead to $p(y|x)$ which may vary largely due to small changes in input
- Encourage weights to be small by adding a penalty
- L2 regularization:

$$\|\vec{\theta}\|^2 = \left(\sqrt{\theta_1^2 + \theta_2^2 + \dots + \theta_K^2} \right)^2 = \sum_{k=1}^K \theta_k^2$$

Loss Function (or Objective Function)

- Goal: Find parameters that minimize this Loss Function,
Negative Log-Likelihood + Regularizer

$$L(\vec{\theta}) = - \sum_{n=1}^N \log(p(y_n | x_n; \vec{\theta})) + \|\vec{\theta}\|^2$$

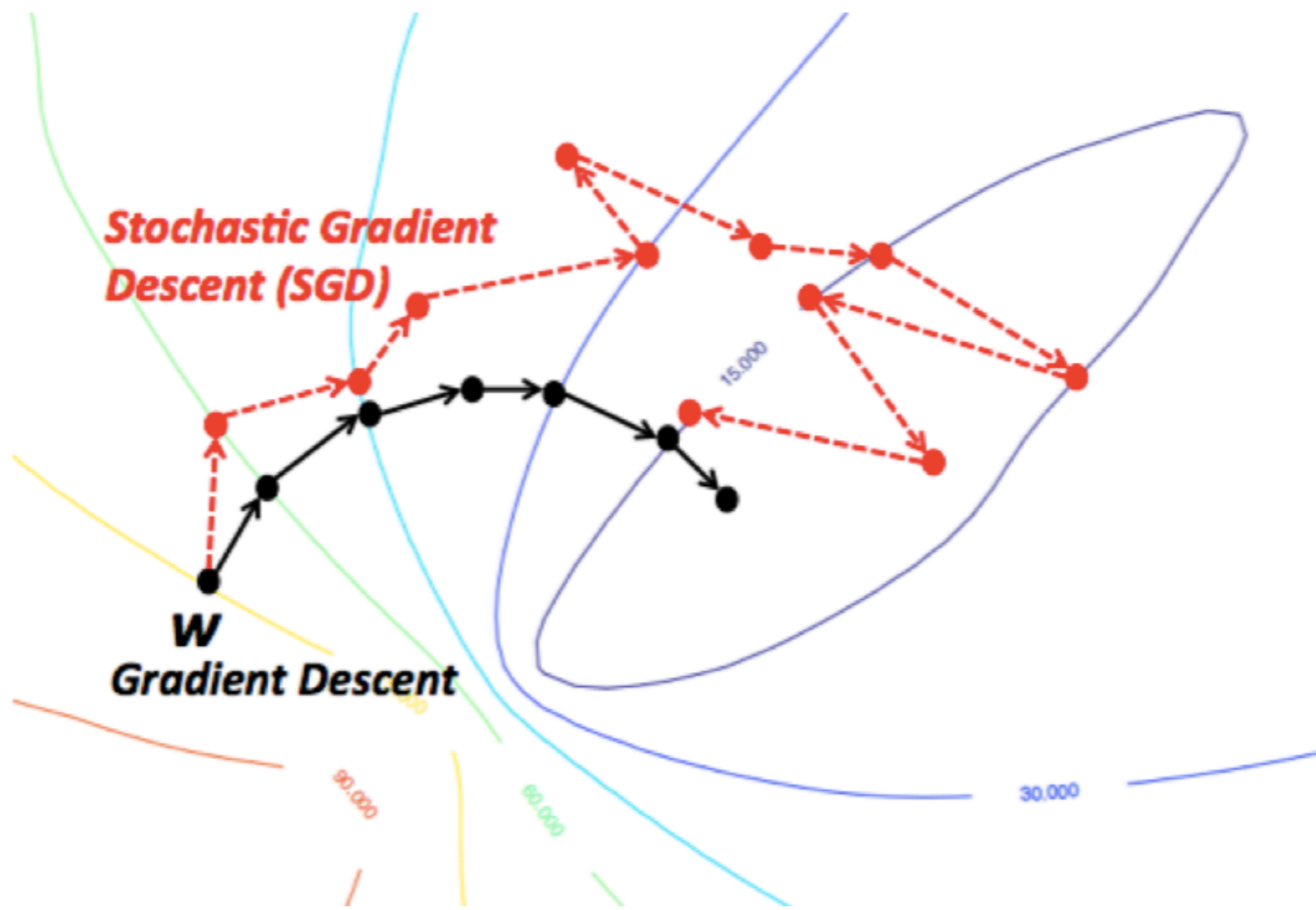
- We can use various optimization techniques. Think back to your Calculus class.

Training by Gradient Descent

- The gradient of a function points to the direction of steepest increase in that function

$$\nabla L(\vec{\theta}) = \left[\frac{\partial L(\vec{\theta})}{\partial \theta_1}; \frac{\partial L(\vec{\theta})}{\partial \theta_2}; \dots; \frac{\partial L(\vec{\theta})}{\partial \theta_K} \right]$$

- Gradient Descent Algorithm: start with some random parameter, keep going in the opposite gradient direction
- Like how you ski down a mountain



Training by Gradient Descent

- Exercise: Compute partial derivatives of loss function

$$L(\vec{\theta}) = - \sum_{n=1}^N \log(p(y_n | x_n; \vec{\theta})) + \|\vec{\theta}\|^2$$

- Regularizer part:

$$\frac{\partial \|\vec{\theta}\|^2}{\partial \theta_1} = \frac{\partial \left(\sqrt{\theta_1^2 + \theta_2^2 + \dots + \theta_K^2} \right)}{\partial \theta_1} = \frac{\partial (\theta_1^2 + \theta_2^2 + \dots + \theta_K^2)}{\partial \theta_1} = \frac{\partial (\theta_1^2)}{\partial \theta_1} = 2\theta_1$$

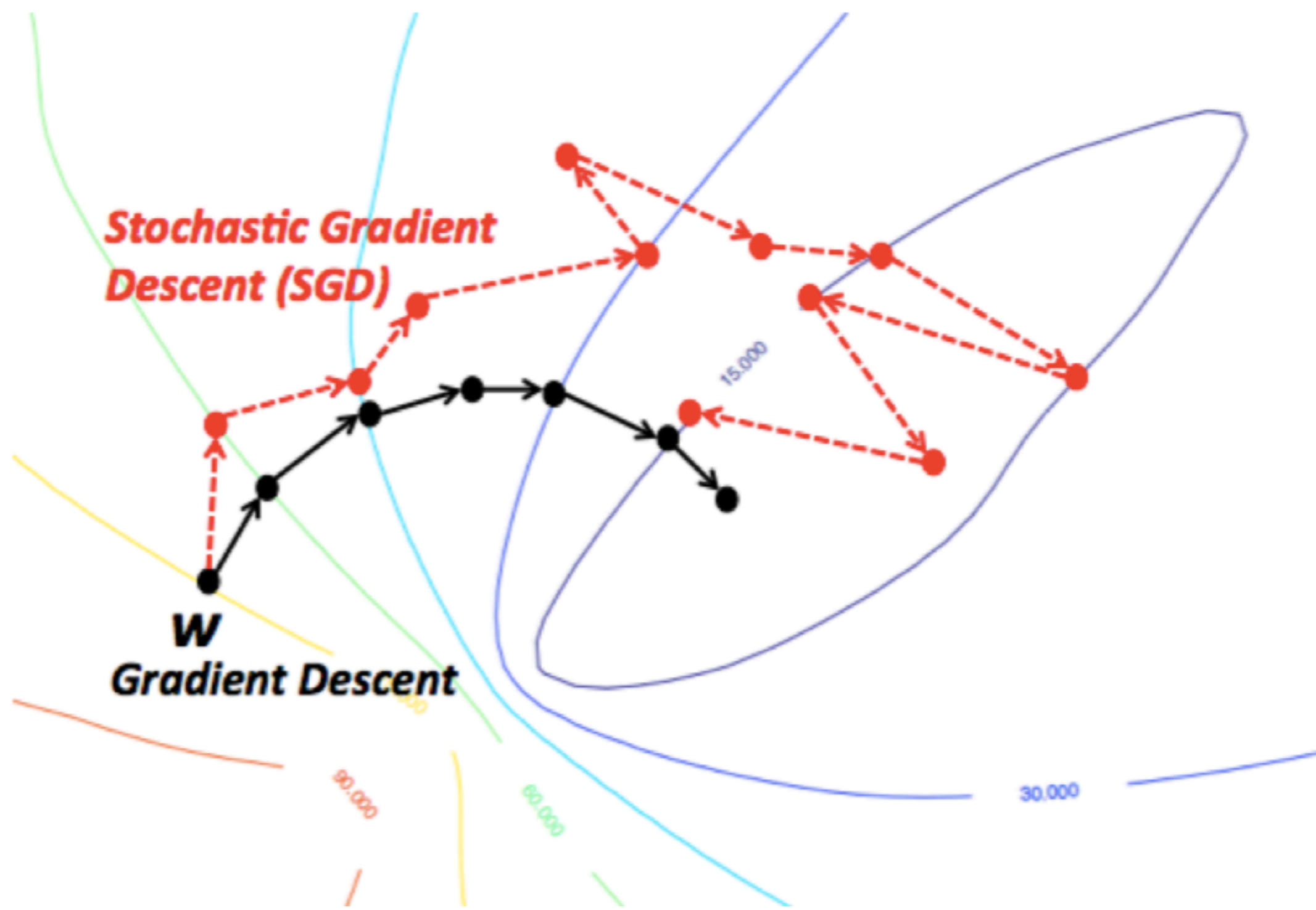
- Likelihood part:

$$\frac{\partial \left(- \sum_{n=1}^N \log p(y_n | x_n; \vec{\theta}) \right)}{\partial \theta_1} = \sum_{n=1}^N \frac{\partial \left(- \log p(y_n | x_n; \vec{\theta}) \right)}{\partial \theta_1}$$

Gradient Descent and Stochastic Gradient Descent (SGD)

- Gradient Descent Algorithm:
 - Start with some parameter $\vec{\theta}$
 - While not converged:
 - Update parameter $\vec{\theta} := \vec{\theta} - \eta \nabla L(\vec{\theta})$
- SGD:

$$\sum_{n=1}^N \frac{\partial \left(-\log p(y_n | x_n; \vec{\theta}) \right)}{\partial \theta_1} \approx \sum_{n:\text{subset}} \frac{\partial \left(-\log p(y_n | x_n; \vec{\theta}) \right)}{\partial \theta_1}$$



Summary

- Log-Linear Model has this form:

$$p(y_n | x_n; \vec{\theta}) = \frac{1}{Z(x)} \exp\left(\sum_{k=1}^K \theta_k f(x_n, y_n)\right)$$

- Features enable us to incorporate knowledge of the problem
- Given training data, we fit parameters to maximize likelihood (optionally with a regularization term)

$$p(\text{training data}; \vec{\theta}) = \prod_{n=1}^N p(y_n | x_n; \vec{\theta})$$

Interactive Visualization

- <http://www.cs.jhu.edu/~jason/465/hw-prob/loglin/#1>

Log-Likelihood Scores

Current LL:  -83.178

Data & Model Options

Change the data

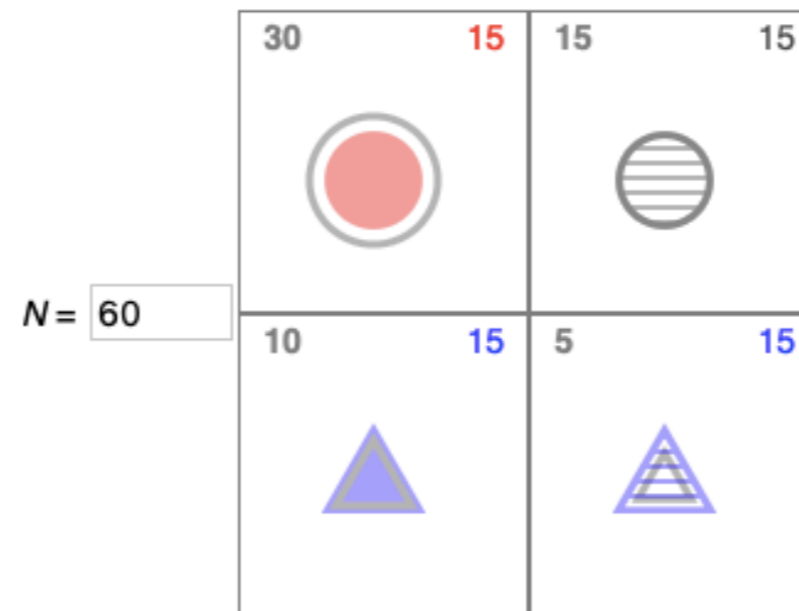
New random challenge

New counts

Regularization

- None
- ℓ_1
- ℓ_2

Type Counts: Observed and Expected



Hints

Show gradient

Feature Weights

circle



0

solid



0

Zero weights