An aerial photograph of a vast, snow-covered mountain range. The peaks are rugged and partially covered in dark, rocky patches. The valleys and slopes are filled with thick, white snow, creating a high-contrast landscape. The sky is a pale, hazy blue, suggesting a clear but slightly overcast day. The overall scene conveys a sense of scale and natural beauty.

S&DS 265 / 565  
Introductory Machine Learning

# Stochastic Gradient Descent

September 19

Yale

# Goings on

- Assn 1 due today at midnight
- Assn 2 will be posted this afternoon
- Quiz 2 next Thursday

# Outline for today

- Continue discussion of generative/discriminative
- Stochastic gradient descent
- Application to logistic regression
- Regularization
- Learning rate and scaling
- Jupyter notebook example

# Stochastic gradient descent

- Suppose that we want to fit a really big model, where the number of samples  $n$  and number of variables  $p$  are very large
- The classical algorithms in standard software packages will fail
- How can we train such models?

# Recall: Classification tasks

- Ad click-through prediction. Predict whether or not a user will click on an ad presented. Used for ranking ads and setting prices.

The screenshot shows a Google search results page for the query "tax preparation services". The page features several sponsored advertisements at the top, including one from "accountant prices.com" and another from "CohnReznick Tax Advisors". Below the ads, there are organic search results for "Businesses" in the area, listing "Cannon John" and "Lao Americas Tax & Financial Agency LLC". On the right side of the page, a map is displayed, showing the location of the search results in the New Haven area, with markers for "Hartford Tax Service" and "Yale University".

# Recall: Classification tasks

- Ad click-through prediction. Predict whether or not a user will click on an ad presented. Used for ranking ads and setting prices.

## Ad targeting

### How ads are targeted to your site



NEXT: ABOUT THE AD AUCTION >

Google automatically delivers ads that are **targeted** to your content or audience. We do this in several ways:

- **Contextual targeting**

Our technology uses such factors as keyword analysis, word frequency, font size, and the overall link structure of the web, in order to determine what a webpage is about and precisely match Google ads to each page.

- **Placement targeting**

With placement targeting, advertisers choose specific **ad placements**, or subsections of publisher websites, on which to run their ads. Ads that are placement-targeted may not be precisely related to the content of a page, but are hand-picked by advertisers who've determined a match between what your users are interested in and what they have to offer.

- **Personalized advertising**

Personalized advertising enables advertisers to reach users based on their interests, demographics (e.g., "sports enthusiasts") and [other criteria](#). To opt out of personalized advertising, users can change their controls in [Ads Settings](#) [↗](#).

- **Language targeting**

Our technology can also determine the primary language of a page. If your content is in a [language supported by our program](#), AdSense will target ads in the appropriate language to your content. We may look at the language of the pages a user is currently viewing, or has recently viewed, to determine which ads to show. In this case, AdSense may target ads in the user's detected language rather than in the language of your content. Learn more about [ad targeting by language](#).

# Example

- We want to classify ads according to whether or not they will be clicked on by a user
- We have a very large collection of training data
- Ads are represented in terms of a sparse list of features

1 | 5 : 1.1789641e-01    39 : 6.0373064e-02    45 : 1.3163488e-01

- The dataset is too large to load into memory, and the number of features is also very large
- New data are continually arriving
- How can we efficiently train a classifier?

# Online learning

We will introduce a method that

- Reads in the data points one (or a few) at a time
- Updates the model for each sample
- Exploits sparsity of the features
- Uses little memory, never reads in the entire dataset



# Stochastic gradient descent

Initialize all parameters to zero:  $\beta_j = 0, j = 1, \dots, p$ .

Read through the data one record at a time, and update the model.

- 1 Read data item  $x$
- 2 Make a prediction  $\hat{y}(x)$
- 3 Observe the true response/label  $y$
- 4 Update the parameters  $\beta$  so  $\hat{y}$  is closer to  $y$

# Stochastic gradient descent

To begin, suppose we are doing *linear regression*. We initialize all parameters to zero:  $\beta_j = 0, j = 1, \dots, p$ .

We read through the data one record at a time, and update the model.

- 1 Read data item  $x$
- 2 Make a prediction  $\hat{y}(x) = \sum_{j=1}^p \beta_j x_j$
- 3 Observe the true response/label  $y$
- 4 Update the parameters  $\beta$  so  $\hat{y}$  is closer to  $y$

# SGD idea

Here's the idea:

- For each parameter  $\beta_j$ , see what happens to the loss if that parameter is increased a little bit.
- If the loss goes down (up), then increase (decrease)  $\beta_j$  proportionately
- Do this simultaneously for all of the parameters
- Rinse and repeat

# SGD idea

Change  $\beta_j$  by a little bit:

$$\beta_j \rightarrow \beta_j + \varepsilon$$

What happens to the squared error?

$$\begin{aligned}(y - \hat{y})^2 &\rightarrow (y - \hat{y} - \varepsilon x_j)^2 \\ &\approx (y - \hat{y})^2 - 2(y - \hat{y})x_j \varepsilon \\ &= (y - \hat{y})^2 + \underbrace{-2(y - \hat{y})x_j}_g \varepsilon\end{aligned}$$

# SGD idea

We then change the parameter as follows:

$$\begin{aligned}\beta_j &\rightarrow \beta_j - \eta g \\ &= \beta_j - \eta \underbrace{(-2(y - \hat{y})x_j)}_g \\ &= \beta_j + \eta 2(y - \hat{y})x_j\end{aligned}$$

with  $\eta$  a small number. So, we are taking

$$\varepsilon = -\eta g$$

# SGD idea

Why is this a good idea? With this choice of  $\varepsilon$  the squared error decreases:

$$\begin{aligned}(y - \hat{y})^2 &\rightarrow (y - \hat{y} - \varepsilon x_j)^2 \\ &\approx (y - \hat{y})^2 - \eta g^2 \\ &< (y - \hat{y})^2\end{aligned}$$

so we're moving “downhill”

# SGD for general loss

Suppose  $L(y, \beta^T x)$  is the loss for an input  $(x, y)$ , e.g.,  $(y - \beta^T x)^2$

SGD update:

$$\beta_j \leftarrow \beta_j - \eta \frac{\partial L(y, \beta^T x)}{\partial \beta_j}$$

$$\boldsymbol{\beta} \leftarrow \boldsymbol{\beta} - \eta \nabla_{\boldsymbol{\beta}} L(y, \boldsymbol{\beta}^T x) \quad (\text{vector notation})$$

- $\eta$  is the *learning rate* or “step size”
- Needs to be chosen carefully, getting smaller over time

# Gradient descent for general loss

If  $L(\beta)$  is the loss function over subset of training set:

$$\begin{aligned}L(\beta + \eta \mathbf{v}) &\approx L(\beta) + \eta \mathbf{v}^T \nabla L(\beta) \\L(\beta - \eta \nabla L(\beta)) &\approx L(\beta) - \eta \|\nabla L(\beta)\|^2\end{aligned}$$

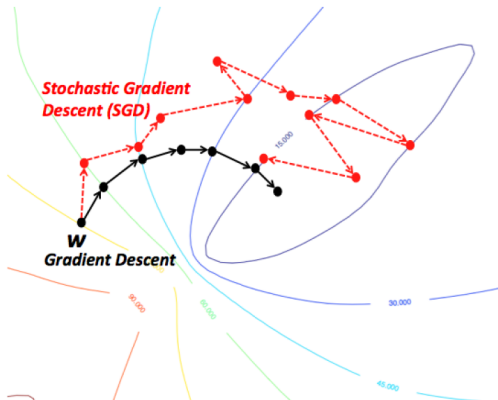
This is why gradient descent is going downhill — if  $\eta$  is small enough.

“Batch” gradient descent uses the entire training set in each step of gradient descent.

*Stochastic* gradient descent computes a quick approximation to this gradient, using only a single or a small “mini-batch” of data points



# Batch vs. stochastic gradient descent



<https://wikidocs.net/3413>

# SGD for logistic regression

SGD Update:

$$\beta_j \longleftarrow \beta_j + \eta(y - p(x))x_j$$

$$\beta_j x_j \longleftarrow \beta_j x_j + \eta(y - p(x))x_j^2$$

$$p(x) = \frac{1}{1 + \exp(-\beta^T x)}$$

Case checking:

- Suppose  $y = 1$  and probability  $p(x)$  is high?

# SGD for logistic regression

SGD Update:

$$\beta_j \longleftarrow \beta_j + \eta(y - p(x))x_j$$

$$\beta_j x_j \longleftarrow \beta_j x_j + \eta(y - p(x))x_j^2$$

$$p(x) = \frac{1}{1 + \exp(-\beta^T x)}$$

Case checking:

- Suppose  $y = 1$  and probability  $p(x)$  is high? *small change*
- Suppose  $y = 1$  and probability  $p(x)$  is small?

# SGD for logistic regression

SGD Update:

$$\beta_j \longleftarrow \beta_j + \eta(y - p(x))x_j$$

$$\beta_j x_j \longleftarrow \beta_j x_j + \eta(y - p(x))x_j^2$$

$$p(x) = \frac{1}{1 + \exp(-\beta^T x)}$$

Case checking:

- Suppose  $y = 1$  and probability  $p(x)$  is high? *small change*
- Suppose  $y = 1$  and probability  $p(x)$  is small? *big change*  $\uparrow$
- Suppose  $y = 0$  and probability  $p(x)$  is small?

# SGD for logistic regression

SGD Update:

$$\beta_j \longleftarrow \beta_j + \eta(y - p(x))x_j$$

$$\beta_j x_j \longleftarrow \beta_j x_j + \eta(y - p(x))x_j^2$$

$$p(x) = \frac{1}{1 + \exp(-\beta^T x)}$$

Case checking:

- Suppose  $y = 1$  and probability  $p(x)$  is high? *small change*
- Suppose  $y = 1$  and probability  $p(x)$  is small? *big change*  $\uparrow$
- Suppose  $y = 0$  and probability  $p(x)$  is small? *small change*
- Suppose  $y = 0$  and probability  $p(x)$  is big?

# SGD for logistic regression

SGD Update:

$$\beta_j \longleftarrow \beta_j + \eta(y - p(x))x_j$$

$$\beta_j x_j \longleftarrow \beta_j x_j + \eta(y - p(x))x_j^2$$

$$p(x) = \frac{1}{1 + \exp(-\beta^T x)}$$

Case checking:

- Suppose  $y = 1$  and probability  $p(x)$  is high? *small change*
- Suppose  $y = 1$  and probability  $p(x)$  is small? *big change*  $\uparrow$
- Suppose  $y = 0$  and probability  $p(x)$  is small? *small change*
- Suppose  $y = 0$  and probability  $p(x)$  is big? *big change*  $\downarrow$

# SGD: choice of step size

In theory, we need to let the step size  $\eta$  decrease as the algorithm progresses.

This prevents the estimates from oscillating back and forth without converging.

# Demo

Open the demo notebook `sgd.ipynb` and follow along...



# SGD: Scaling

We generally want to “standardize” each variable — subtract out the mean and divide by the standard deviation

$$x_j \leftarrow \frac{x_j - \text{mean}(x_j)}{\sqrt{\text{var}(x_j)}}$$

But this involves “looking ahead” to compute the mean and variance, and destroys the online property of the algorithm

# SGD: Scaling

We generally want to “standardize” each variable — subtract out the mean and divide by the standard deviation

$$x_j \leftarrow \frac{x_j - \text{mean}(x_j)}{\sqrt{\text{var}(x_j)}}$$

But this involves “looking ahead” to compute the mean and variance, and destroys the online property of the algorithm

Solution: The mean and variance can be updated in an online manner, in constant time, by storing auxiliary variables for each component  $j$ .

# SGD: Regularization

A “ridge” penalty  $\lambda \sum_{j=1}^p \beta_j^2$  is easily handled.

Gradient changes by an additive term  $2\lambda\beta$ . Update becomes

$$\begin{aligned}\beta_j &\longleftarrow \beta_j + \eta\{(y - p(x))x_j - 2\lambda\beta_j\} \\ &= (1 - 2\eta\lambda)\beta_j + \eta(y - p(x))x_j\end{aligned}$$

Observe that this “does the right thing” whether  $\beta_j$  wants to be large positive or negative.

- *The penalty shrinks  $\beta_j$  toward zero*

# What did we learn today?

- Stochastic gradient descent is a simple algorithm that can be applied to large classification and regression problems
- A parameter is updated according to how much the loss changes when that parameter is changed by a little bit
- This is the “go to” algorithm for fitting large or complex machine learning models
- Choosing the learning rate is a little tricky