

EE4077 Fundamentals of Machine Learning

Logistic Regression

Fall 2021

EE–Marmara University

Outline

1. Logistic Regression Model

2. Loss Function and Parameter Estimation

3. Gradient Descent

Logistic Regression Model

How does naive Bayes work?

Examine the classification rule for naive Bayes

$$y^* = \arg \max_c \left(\log \pi_c + \sum_k x_k \log \theta_{ck} \right)$$

For binary classification, we thus determine the label based on the sign of

$$\log \pi_1 + \sum_k x_k \log \theta_{1k} - \left(\log \pi_2 + \sum_k x_k \log \theta_{2k} \right)$$

This is just a linear function of the features (word-counts) $\{x_k\}$

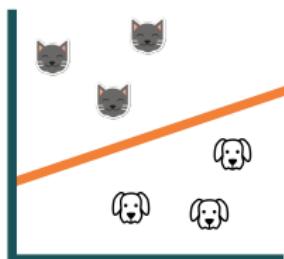
$$w_0 + \sum_k x_k w_k$$

where we “absorb” $w_0 = \log \pi_1 - \log \pi_2$ and $w_k = \log \theta_{1k} - \log \theta_{2k}$.

Intuition: Logistic Regression

Learn the equation of the decision boundary $\mathbf{w}^\top \mathbf{x} = 0$ such that

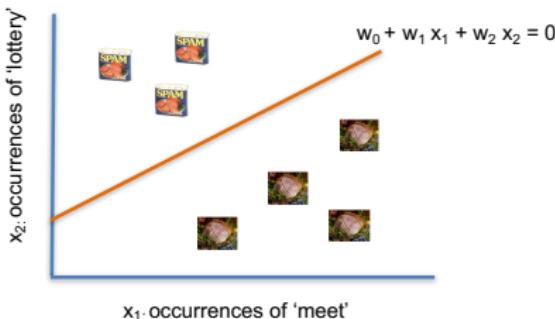
- If $\mathbf{w}^\top \mathbf{x} \geq 0$ declare $y = 1$ (cat)
- If $\mathbf{w}^\top \mathbf{x} < 0$ declare $y = 0$ (dog)



$y = 0$ for dog, $y = 1$ for cat

Back to spam vs. ham classification...

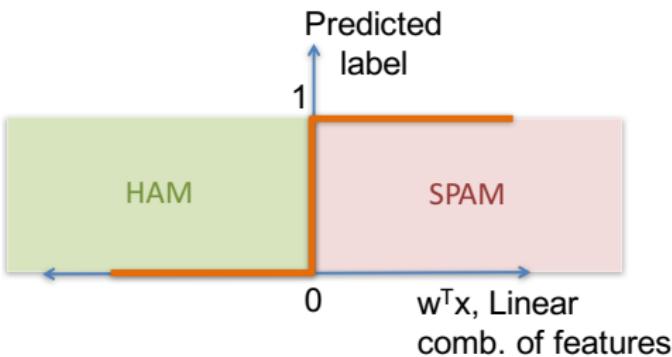
- $x_1 = \#$ of times 'meet' appears in an email
- $x_2 = \#$ of times 'lottery' appears in an email
- Define feature vector $\mathbf{x} = [1, x_1, x_2]$
- Learn the decision boundary $w_0 + w_1 x_1 + w_2 x_2 = 0$ such that
 - If $\mathbf{w}^\top \mathbf{x} \geq 0$ declare $y = 1$ (spam)
 - If $\mathbf{w}^\top \mathbf{x} < 0$ declare $y = 0$ (ham)



Key Idea: If 'meet' appears few times and 'lottery' appears many times than the email is spam

Visualizing a linear classifier

- $x_1 = \#$ of times 'lottery' appears in an email
- $x_2 = \#$ of times 'meet' appears in an email
- Define feature vector $\mathbf{x} = [1, x_1, x_2]$
- Learn the decision boundary $w_0 + w_1x_1 + w_2x_2 = 0$ such that
 - If $\mathbf{w}^\top \mathbf{x} \geq 0$ declare $y = 1$ (spam)
 - If $\mathbf{w}^\top \mathbf{x} < 0$ declare $y = 0$ (ham)



$$y = 1 \text{ for spam, } y = 0 \text{ for ham}$$

Your turn

Suppose you see the following email:

CONGRATULATIONS!! Your email address have won you the lottery sum of US\$2,500,000.00 USD to claim your prize, contact your office agent (Arthur walter) via email claims2155@yahoo.com.hk or call +44 704 575 1113

Keywords are [lottery, prize, office, email]

The given weight vector is $\mathbf{w} = [0.3, 0.3, -0.1, -0.04]^\top$

Will we predict that the email is spam or ham?

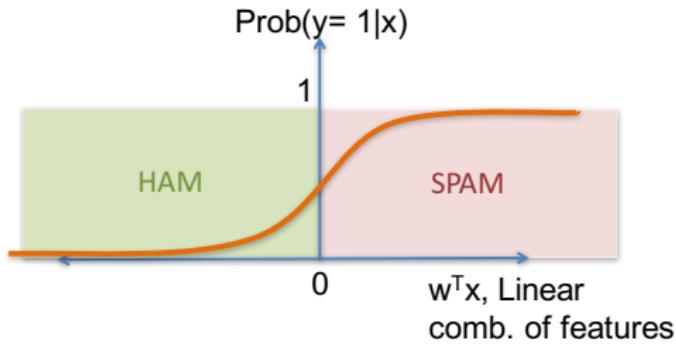
$$\mathbf{x} = [1, 1, 1, 2]^\top$$

$$\mathbf{w}^\top \mathbf{x} = 0.3 * 1 + 0.3 * 1 - 0.1 * 1 - 0.04 * 2 = 0.42 > 0$$

so we predict spam!

Intuition: Logistic Regression

- Suppose we want to output the **probability** of an email being spam/ham instead of just 0 or 1
- This gives information about the confidence in the decision
- Use a function $\sigma(\mathbf{w}^\top \mathbf{x})$ that maps $\mathbf{w}^\top \mathbf{x}$ to a value between 0 and 1



Probability that predicted label is 1 (spam)

Key Problem: Finding optimal weights \mathbf{w} that accurately predict this probability for a new email

Formal Setup: Binary Logistic Classification

- Input/features: $\mathbf{x} = [1, x_1, x_2, \dots, x_D] \in \mathbb{R}^{D+1}$
- Output: $y \in \{0, 1\}$
- Training data: $\mathcal{D} = \{(\mathbf{x}_n, y_n), n = 1, 2, \dots, N\}$
- Model:

$$p(y = 1 | \mathbf{x}; \mathbf{w}) = \sigma[g(\mathbf{x})]$$

where

$$g(\mathbf{x}) = w_0 + \sum_d w_d x_d = \mathbf{w}^\top \mathbf{x}$$

and $\sigma[\cdot]$ stands for the **sigmoid** function

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$

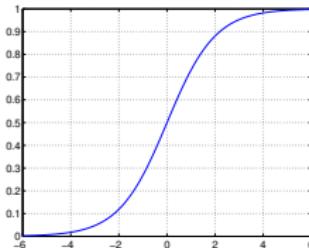
Why the sigmoid function?

What does it look like?

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$

where

$$a = \mathbf{w}^\top \mathbf{x}$$

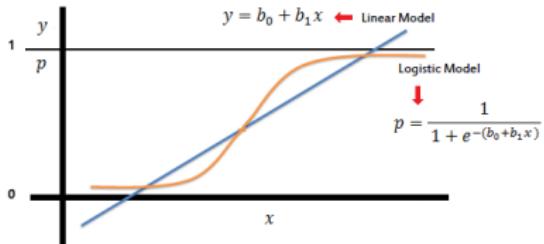


Sigmoid properties

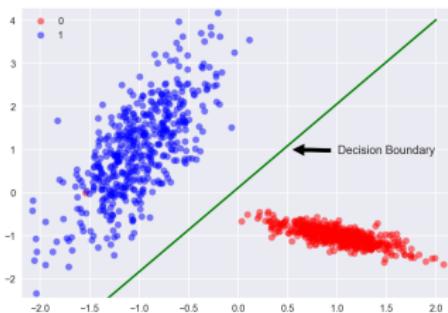
- Bounded between 0 and 1 ← thus, interpretable as probability
- Monotonically increasing ← thus, usable to derive classification rules
 - $\sigma(a) \geq 0.5$, positive (classify as '1')
 - $\sigma(a) < 0.5$, negative (classify as '0')
- Nice computational properties ← as we will see soon

Comparison to Linear Regression

Sigmoid function returns values in $[0,1]$



Decision boundary is linear



Your turn

Suppose you see the following email:

CONGRATULATIONS!! Your email address have won you the lottery sum of US\$2,500,000.00 USD to claim your prize, contact your office agent (Athur walter) via email claims2155@yahoo.com.hk or call +44 704 575 1113

Keywords are [lottery, prize, office, email]

The given weight vector is $\mathbf{w} = [0.3, 0.3, -0.1, -0.04]^\top$

What is the probability that the email is spam?

$$\mathbf{x} = [1, 1, 1, 2]^\top$$

$$\mathbf{w}^\top \mathbf{x} = 0.3 * 1 + 0.3 * 1 - 0.1 * 1 - 0.04 * 2 = 0.42 > 0$$

$$\Pr(y = 1 | \mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + e^{-0.42}} = 0.603$$

Loss Function and Parameter Estimation

How do we optimize the weight vector w?

Learn from experience

- get a lot of spams
- get a lot of hams

But what to optimize?



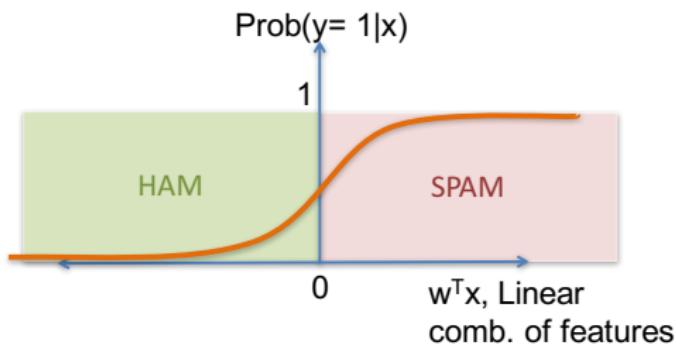
Likelihood function

Probability of a single training sample (\mathbf{x}_n, y_n) ...

$$p(y_n | \mathbf{x}_n; \mathbf{w}) = \begin{cases} \sigma(\mathbf{w}^\top \mathbf{x}_n) & \text{if } y_n = 1 \\ 1 - \sigma(\mathbf{w}^\top \mathbf{x}_n) & \text{otherwise} \end{cases}$$

Simplify, using the fact that y_n is either 1 or 0

$$p(y_n | \mathbf{x}_n; \mathbf{w}) = \sigma(\mathbf{w}^\top \mathbf{x}_n)^{y_n} [1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)]^{1-y_n}$$



Probability that predicted label is 1 (spam)

Log Likelihood or Cross Entropy Error

Log-likelihood of the whole training data \mathcal{D}

$$P(\mathcal{D}) = \prod_{n=1}^N p(y_n | \mathbf{x}_n; \mathbf{w}) = \prod_{n=1}^N \{\sigma(\mathbf{w}^\top \mathbf{x}_n)^{y_n} [1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)]^{1-y_n}\}$$
$$\log P(\mathcal{D}) = \sum_n \{y_n \log \sigma(\mathbf{w}^\top \mathbf{x}_n) + (1 - y_n) \log [1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)]\}$$

It is convenient to work with its negation, which is called the
cross-entropy error function

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^\top \mathbf{x}_n) + (1 - y_n) \log [1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)]\}$$

How to find the optimal parameters for logistic regression?

We will minimize the error function

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^\top \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)]\}$$

However, this function is complex and we cannot find the simple solution as we did in Naive Bayes. So we need to use **numerical** methods.

- Numerical methods are messier, in contrast to cleaner closed-form solutions.
- In practice, we often have to tune a few optimization parameters — patience is necessary.
- A popular method: **gradient descent** and its variants.

Gradient Descent

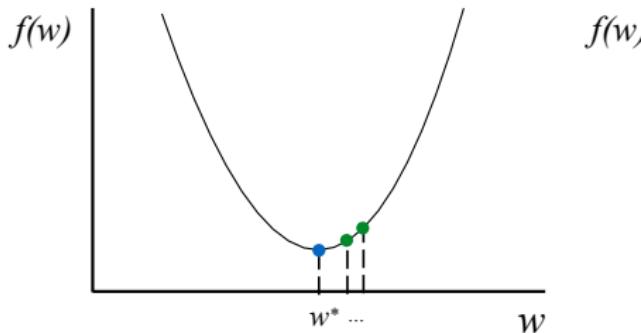
Gradient descent algorithm

Start at a random point.

Repeat:

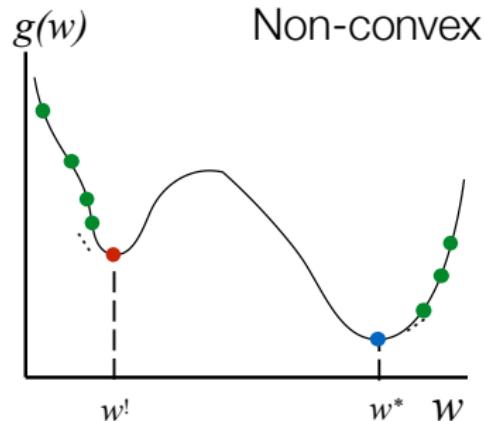
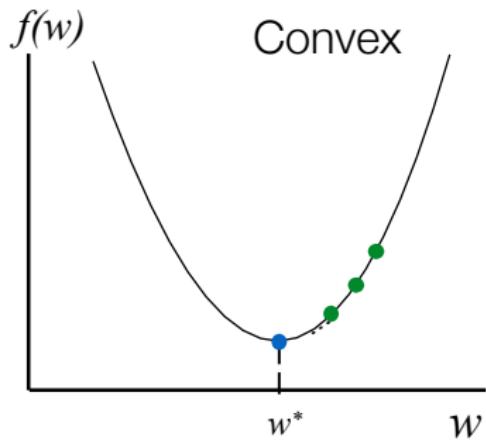
- Determine a descent direction.
- Choose a step size.
- Update.

Until stopping criterion is reached.



Can we converge to the optimum?

Gradient descent (with proper step size) converges to the global optimum for when **minimizing a convex function**.



Any local minimum is also a global minimum.

Multiple local minima may exist.

Linear regression, ridge regression, and logistic regression are all convex!

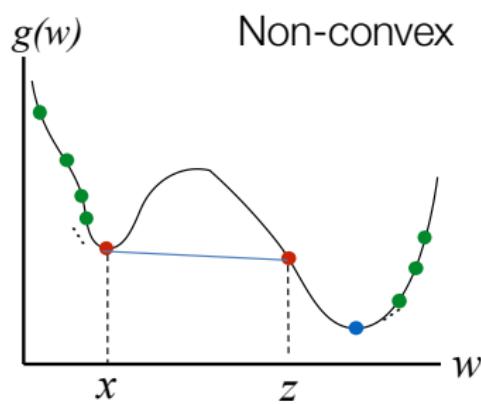
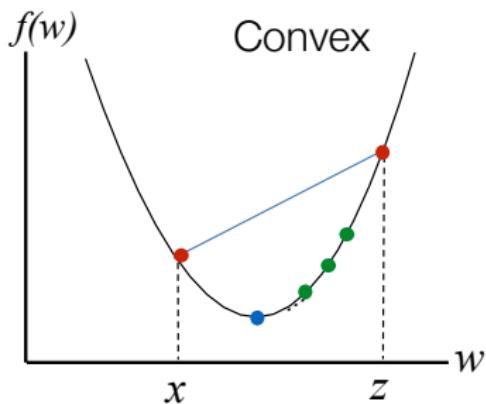
Convexity

A function $f : \mathbb{R}^k \rightarrow \mathbb{R}$ is **convex** if for any $x, z \in \mathbb{R}^k$ and $t \in [0, 1]$,

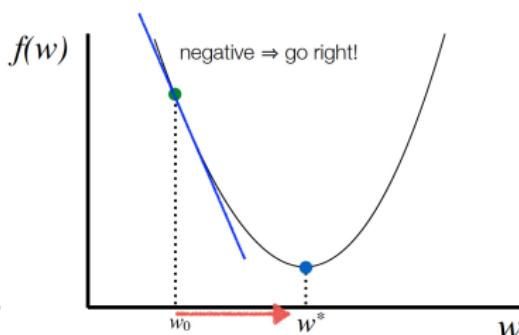
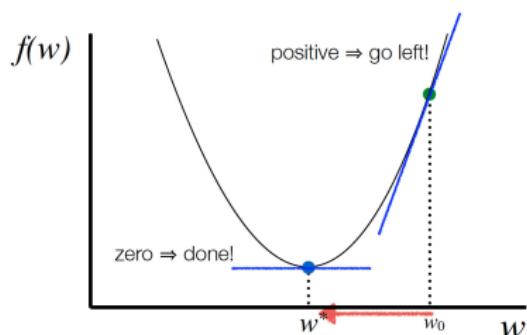
$$\underbrace{f(tx + (1-t)z)}_{\text{Function value at a point between } x \text{ and } z.} \leq \underbrace{tf(x) + (1-t)f(z)}_{\text{Line drawn between } f(x) \text{ and } f(z)}. \quad .$$

Function value at a point between x and z . Line drawn between $f(x)$ and $f(z)$

- f always lies below a line drawn between two of its points.
- If it exists, the Hessian $\frac{d^2f}{dx^2}$ is a positive-(semi)definite matrix.



Why do we move in the direction opposite the gradient?



We can only move in two directions
Negative slope is direction of descent!

Step Size

Update Rule: $w_{i+1} = w_i - \alpha_i \frac{df}{dw}(w_i)$

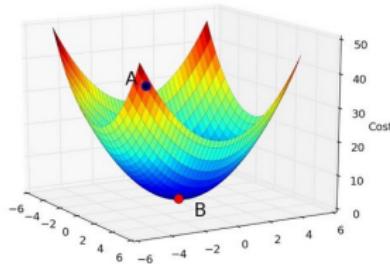
Negative Slope

We have seen this before!

(Batch) gradient descent for linear regression

$$RSS(\mathbf{w}) = \sum_n [y_n - \mathbf{w}^\top \mathbf{x}_n]^2 = \left\{ \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - 2 (\mathbf{X}^\top \mathbf{y})^\top \mathbf{w} \right\} + \text{const}$$

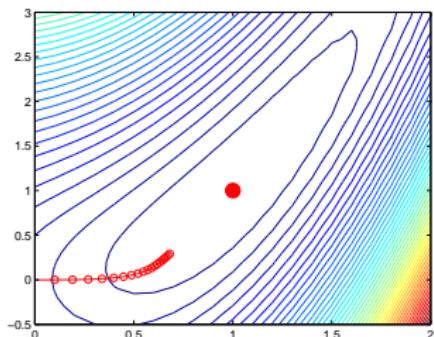
- Initialize \mathbf{w} to $\mathbf{w}^{(0)}$ (e.g., randomly);
set $t = 0$; choose $\eta > 0$
- Loop *until convergence*
 1. Compute the gradient
$$\nabla RSS(\mathbf{w}) = \mathbf{X}^\top (\mathbf{X} \mathbf{w}^{(t)} - \mathbf{y})$$
 2. Update the parameters
$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla RSS(\mathbf{w})$$
 3. $t \leftarrow t + 1$



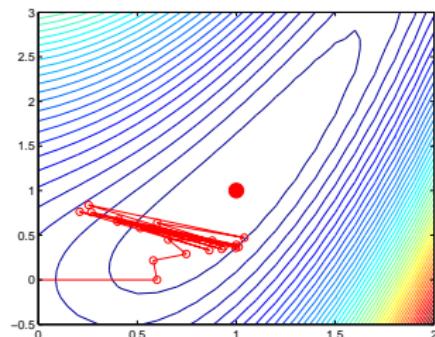
Impact of step size

Choosing the right η is important

small η is too slow?

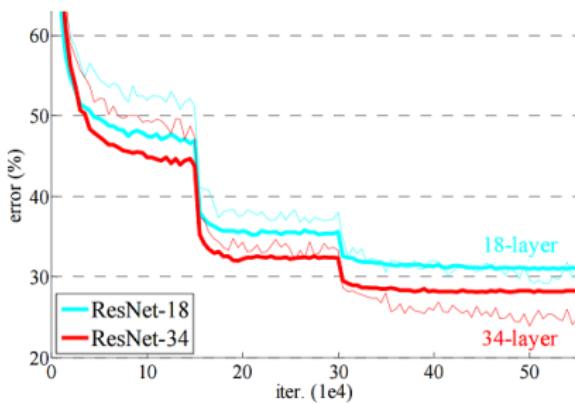


large η is too unstable?



How to choose η in practice?

- Try 0.0001, 0.001, 0.01, 0.1 etc. on a validation dataset and choose the one that gives fastest, stable convergence
- Reduce η by a constant factor (eg. 10) when learning saturates so that we can reach closer to the true minimum.
- More advanced learning rate schedules such as AdaGrad, Adam, AdaDelta are used in practice.



Gradient descent for a general function

General form for minimizing $f(\theta)$

$$\theta^{t+1} \leftarrow \theta^t - \eta \frac{\partial f}{\partial \theta} \Big|_{\theta=\theta^t}$$

- η is **step size**, also called the **learning rate** – how far we go in the direction of the negative gradient
 - Step size needs to be chosen carefully to ensure convergence.
 - Step size can be adaptive, e.g., we can use **line search**
- We are **minimizing** a function, hence the subtraction $(-\eta)$
- With a **suitable** choice of η , we converge to a stationary point

$$\frac{\partial f}{\partial \theta} = 0$$

- Stationary point not always global minimum (but happy when convex)
- Popular variant called **stochastic** gradient descent

Gradient descent update for Logistic Regression

Finding the gradient of $\mathcal{E}(\mathbf{w})$ looks very hard, but it turns out to be simple and intuitive.

Let's start with the derivative of the sigmoid function $\sigma(a)$:

$$\begin{aligned}\frac{d}{da}\sigma(a) &= \frac{d}{da}(1 + e^{-a})^{-1} \\&= \frac{-1}{(1 + e^{-a})^2} \frac{d}{da}(1 + e^{-a}) \\&= \frac{e^{-a}}{(1 + e^{-a})^2} \\&= \frac{1}{1 + e^{-a}} \frac{e^{-a}}{1 + e^{-a}} \\&= \frac{1}{1 + e^{-a}} \frac{1 + e^{-a} - 1}{1 + e^{-a}} \\&= \sigma(a)[1 - \sigma(a)]\end{aligned}$$

Gradients of the cross-entropy error function

Cross-entropy Error Function

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^\top \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)]\}$$

$$\frac{d}{da} \sigma(a) = \sigma(a)[1 - \sigma(a)]$$

Computing the gradient

$$\begin{aligned}\frac{\partial \mathcal{E}(\mathbf{w})}{\partial \mathbf{w}} &= - \sum_n \left\{ y_n \frac{\sigma(\mathbf{w}^\top \mathbf{x}_n)[1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)]}{\sigma(\mathbf{w}^\top \mathbf{x}_n)} \mathbf{x}_n \right. \\ &\quad \left. - (1 - y_n) \frac{\sigma(\mathbf{w}^\top \mathbf{x}_n)[1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)]}{1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)} \mathbf{x}_n \right\} \\ &= - \sum_n \{y_n [1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)] \mathbf{x}_n - (1 - y_n) \sigma(\mathbf{w}^\top \mathbf{x}_n) \mathbf{x}_n\} \\ &= \sum_n \underbrace{\{\sigma(\mathbf{w}^\top \mathbf{x}_n) - y_n\}}_{\text{Error of the } n\text{th training sample.}} \mathbf{x}_n\end{aligned}$$

Gradient descent for logistic regression

- Choose a proper step size $\eta > 0$
- Iteratively update the parameters following the negative gradient to minimize the error function

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \sum_n \left\{ \sigma(\mathbf{w}^{(t)\top} \mathbf{x}_n) - y_n \right\} \mathbf{x}_n$$

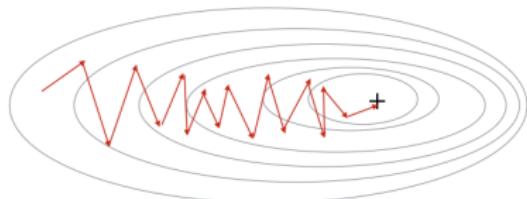
Stochastic gradient descent for logistic regression

- Choose a proper step size $\eta > 0$
- Draw a sample n uniformly at random
- Iteratively update the parameters following the negative gradient to minimize the error function

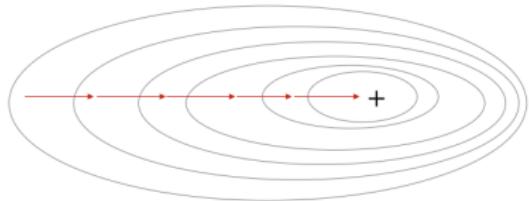
$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \left\{ \sigma(\mathbf{w}^{(t)\top} \mathbf{x}_n) - y_n \right\} \mathbf{x}_n$$

SGD versus Batch GD

Stochastic Gradient Descent



Gradient Descent



- SGD reduces per-iteration complexity since it considers fewer samples.
- But it is noisier and can take longer to converge.

Example: Spam Classification

	free	bank	meet	time	y
Email 1	5	3	1	1	Spam
Email 2	4	2	1	1	Spam
Email 3	2	1	2	3	Ham
Email 4	1	2	3	2	Ham

Perform gradient descent to learn weights \mathbf{w}

- Feature vector for email 1: $\mathbf{x}_1 = [1, 5, 3, 1, 1]^\top$
- Let $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4]$, the matrix of all feature vectors.
- Initial weights $\mathbf{w} = [0.5, 0.5, 0.5, 0.5, 0.5]^\top$
- Prediction

$$[\sigma(\mathbf{w}^\top \mathbf{x}_1), \sigma(\mathbf{w}^\top \mathbf{x}_2), \sigma(\mathbf{w}^\top \mathbf{x}_3), \sigma(\mathbf{w}^\top \mathbf{x}_4)]^\top = [0.996, 0.989, 0.989, 0.989]^\top$$

which can be obtained by computing $\mathbf{w}^\top \mathbf{X}$ and then apply $\sigma(\cdot)$ entrywise, which we abuse the notation and write $\sigma(\mathbf{X}^\top \mathbf{w})$.

Example: Spam Classification, Batch Gradient Descent

	free	bank	meet	time	y
Email 1	5	3	1	1	Spam
Email 2	4	2	1	1	Spam
Email 3	2	1	2	3	Ham
Email 4	1	2	3	2	Ham

Perform gradient descent to learn weights \mathbf{w}

- Prediction $\sigma(\mathbf{X}^\top \mathbf{w}) = [0.996, 0.989, 0.989, 0.989]^\top$
- Difference from labels $\mathbf{y} = [1, 1, 0, 0]^\top$ is

$$\sigma(\mathbf{X}^\top \mathbf{w}) - \mathbf{y} = [-0.004, -0.011, 0.989, 0.989]^\top$$

- Gradient of the first email,

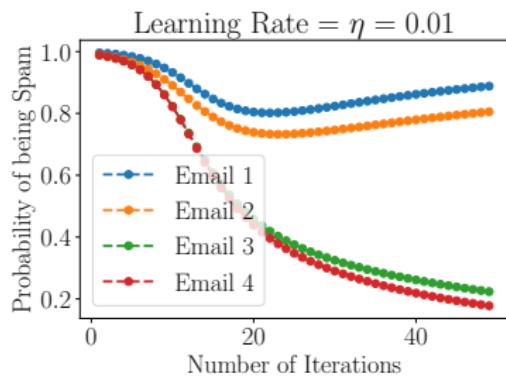
$$\mathbf{g}_1 = (\sigma(\mathbf{w}^\top \mathbf{x}_1) - y_1) \mathbf{x}_1 = -0.004[1, 5, 3, 1, 1]^\top$$

- $\mathbf{w} \leftarrow \mathbf{w} - \underbrace{0.01}_{\text{learning rate}} \sum_n \mathbf{g}_n = \mathbf{w} - \eta \mathbf{X}(\sigma(\mathbf{X}^\top \mathbf{w}) - \mathbf{y})$

notice the similarity with linear regression

Example: Spam Classification, Batch Gradient Descent

	free	bank	meet	time	y
Email 1	5	3	1	1	Spam
Email 2	4	2	1	1	Spam
Email 3	2	1	2	3	Ham
Email 4	1	2	3	2	Ham



Predictions for Emails 3 and 4 are initially close to 1 (spam), but they converge towards the correct value 0 (ham)

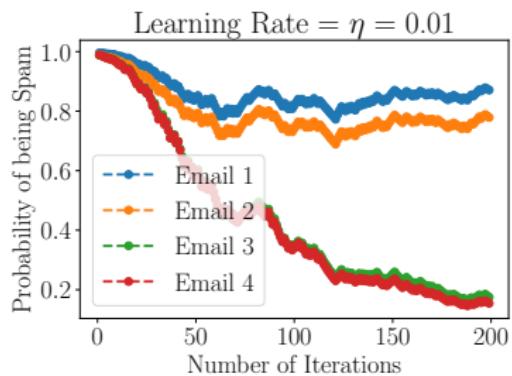
Example: Spam Classification, Stochastic Gradient Descent

	free	bank	meet	time	y
Email 1	5	3	1	1	Spam
Email 2	4	2	1	1	Spam
Email 3	2	1	2	3	Ham
Email 4	1	2	3	2	Ham

- Prediction $\sigma(\mathbf{w}^\top \mathbf{x}_r) = 0.996$ for a randomly chosen email r
- Difference from label $y = 1$ is -0.004
- Gradient is $\mathbf{g}_r = (\sigma(\mathbf{w}^\top \mathbf{x}_r) - y)\mathbf{x}_r = -0.004\mathbf{x}_r$
- $\mathbf{w} \leftarrow \mathbf{w} - 0.01\mathbf{g}_r$

Example: Spam Classification, Stochastic Gradient Descent

	free	bank	meet	time	y
Email 1	5	3	1	1	Spam
Email 2	4	2	1	1	Spam
Email 3	2	1	2	3	Ham
Email 4	1	2	3	2	Ham



Predictions for Emails 3 and 4 are initially close to 1 (spam), but they converge towards the correct value 0 (ham)

Example: Spam Classification, Test Phase

	free	bank	meet	time	y
Email 1	5	3	1	1	Spam
Email 2	4	2	1	1	Spam
Email 3	2	1	2	3	Ham
Email 4	1	2	3	2	Ham

- Final $\mathbf{w} = [0.187, 0.482, 0.179, -0.512, -0.524]^\top$ after 50 batch gradient descent iterations.
- Given a new email with feature vector $\mathbf{x} = [1, 1, 3, 4, 2]$, the probability of the email being spam is estimated as $\sigma(\mathbf{w}^\top \mathbf{x}) = \sigma(-1.889) = 0.13$.
- Since this is less than 0.5 we predict ham.

Contrast Naive Bayes and Logistic Regression

Both classification models are linear functions of features

Joint vs. conditional distribution

Naive Bayes models the **joint** distribution: $P(X, Y) = P(Y)P(X|Y)$

Logistic regression models the **conditional** distribution: $P(Y|X)$

Correlated vs. independent features

Naive Bayes assumes independence of features and multiple occurrences

Logistic Regression implicitly captures correlations when training weights

Generative vs. Discriminative

NB is a **generative** model, LR is a **discriminative** model

Generative Model v.s. Discriminative Model

$\{x : P(Y = 1|X = x) = P(Y = 0|X = x)\}$ is called the **decision boundary** of our data.

Generative classifiers

Model the class-conditional densities $P(Y|X = x)$ explicitly:

$$P(Y = 1|X = x) = \frac{P(X = x|Y = 1)P(Y = 1)}{P(X = x|Y = 1)P(Y = 1) + P(X = x|Y = 0)P(Y = 0)}$$

This means we need to separately estimate both $P(X|Y)$ and $P(Y)$.

Discriminative classifier

Directly model the decision boundary and avoid estimating the conditional probabilities.

Summary

Setup for binary classification

- Logistic Regression models conditional distribution as:
 $p(y = 1|\mathbf{x}; \mathbf{w}) = \sigma[g(\mathbf{x})]$ where $g(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$
- Linear decision boundary: $g(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} = 0$

Minimizing cross-entropy error (negative log-likelihood)

- $\mathcal{E}(b, \mathbf{w}) = -\sum_n \{y_n \log \sigma(b + \mathbf{w}^\top \mathbf{x}_n) + (1 - y_n) \log [1 - \sigma(b + \mathbf{w}^\top \mathbf{x}_n)]\}$
- No closed form solution; must rely on iterative solvers

Numerical optimization

- Gradient descent: simple, scalable to large-scale problems
 - Move in direction opposite of gradient!
 - Gradient of the cross-entropy error takes nice form

What's next?

What about when we want to predict multiple classes?

- Dog vs. cat. vs crocodile
- Movie genres (action, horror, comedy, ...)
- Yelp ratings (1, 2, 3, 4, 5)
- Part of speech tagging (verb, noun, adjective, ...)
- ...