

## CSCI567 Machine Learning (Spring 2019)

Prof. Victor Adamchik

U of Southern California

Jan. 29, 2019

()

January 29, 2019 1 / 50

### Outline

- 1 Gradient Descent
- 2 Linear Classifier and Surrogate Losses
- 3 Perceptron
- 4 Logistic Regression

()

January 29, 2019 3 / 50

## Administration

- TA-1 due this Friday
- PA-2 is released
- Follow Piazza for clarifications
- Notes on grading prog. assignments:
  - ▶ The grading scripts are **NOT** debugging scripts
  - ▶ The scripts provide **enough** feedback information to help you with fixing errors
  - ▶ Do not ask or try to print our test data
  - ▶ You have an unlimited number of submissions - do not abuse the system.

()

January 29, 2019 2 / 50

### Outline

- 1 Gradient Descent
- 2 Linear Classifier and Surrogate Losses
- 3 Perceptron
- 4 Logistic Regression

()

January 29, 2019 4 / 50

## Predicting a continuous outcome variable using past observations

### Key difference from classification

- continuous vs discrete
- measure *prediction errors* differently.
- lead to quite different learning algorithms.

**Linear Regression:** regression with *linear models*:  $f(\mathbf{w}) = \mathbf{w}^T \mathbf{x} = \mathbf{x}^T \mathbf{w}$

# Gradient

The gradient vector  $\nabla f$  points in the direction of greatest rate of increase of  $f$  at a given point.

The the rates of change of  $f$  in all directions is given by

$$\nabla f \cdot \mathbf{u} = \|\nabla f\| \|\mathbf{u}\| \cos \alpha$$

Hence, the direction of *greatest decrease* of  $f$  is the direction opposite to the gradient vector, when  $\alpha = \pi$

We will minimize  $RSS(\mathbf{w})$  using a gradient descent method.

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{n=1}^N (\mathbf{x}_n^T \mathbf{w} - y_n)^2 = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

Three approaches to find the minimum:

- Closed Form (setting gradient to zero)  $\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
- Gradient Descent (GD)
- Stochastic Gradient Descent (SGD)

# Gradient Descent (GD)

**Goal:** minimize  $f(\mathbf{w})$

Consider the definition

$$f'(\mathbf{w}) = \lim_{\Delta x \rightarrow 0} \frac{f(\mathbf{w} + \Delta x) - f(\mathbf{w})}{\Delta x}$$

Our gradient is an estimation of the derivative

$$\nabla f(\mathbf{w}) = \frac{f(\mathbf{w} + \Delta x) - f(\mathbf{w})}{\Delta x}$$

Then we need to move in its *opposite* direction to climb down the function.

$$f(\mathbf{w} + \Delta x) = f(\mathbf{w}) - \Delta x \nabla f(\mathbf{w})$$

## Algorithm: Gradient Descent

**Goal:** minimize  $F(\mathbf{w})$

**Algorithm:** move a bit in the *negative gradient direction*

**initialize**  $\mathbf{w}^{(0)}$

**while not converged do**

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \lambda \nabla F(\mathbf{w}^{(t)})$$

where  $\lambda > 0$  is called *step size* or *learning rate*

- in theory  $\lambda$  should be set in terms of some parameters of  $F$
- in practice we just try several small values
- there are many possible ways to detect convergence.

()

January 29, 2019 9 / 50

## An example

Example:  $F(\mathbf{w}) = 0.5(w_1^2 - w_2)^2 + 0.5(w_1 - 1)^2$ . Gradient is

$$\frac{\partial F}{\partial w_1} = 2(w_1^2 - w_2)w_1 + w_1 - 1 \quad \frac{\partial F}{\partial w_2} = -(w_1^2 - w_2)$$

GD:

- Initialize  $w_1^{(0)}$  and  $w_2^{(0)}$  (to be 0 or *randomly*),  $t = 0$
- do

$$w_1^{(t+1)} \leftarrow w_1^{(t)} - \lambda \left[ 2(w_1^{(t)^2} - w_2^{(t)})w_1^{(t)} + w_1^{(t)} - 1 \right]$$

$$w_2^{(t+1)} \leftarrow w_2^{(t)} - \lambda \left[ -(w_1^{(t)^2} - w_2^{(t)}) \right]$$

$$t \leftarrow t + 1$$

- until  $F(\mathbf{w}^{(t)})$  **does not change much**

()

January 29, 2019 10 / 50

## Why GD?

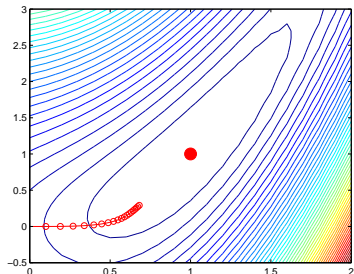
Using the first-order approximation

$$f(w + \Delta x) = f(w) + \Delta x \nabla f(w)$$

we move a bit in the negative gradient direction  $\Delta x = -\lambda \nabla f(w)$

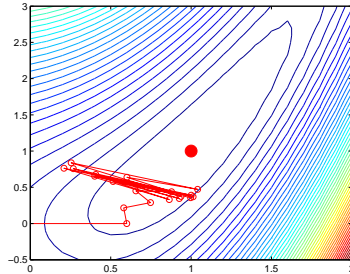
This ensures

$$f(w - \lambda \nabla f(w)) = f(w) - \lambda (\nabla f(w))^2 \leq f(w)$$



$\lambda$  decreases function value

()



reasonable

large  $\lambda$  is unstable

but

January 29, 2019 11 / 50

## Applying GD to Linear Regression

In the previous discussion, we have computed:

$$\frac{1}{2} \nabla RSS(\mathbf{w}) = \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} = \sum_n \mathbf{x}_n (\mathbf{x}_n^T \mathbf{w} - y_n) = \sum_n \mathbf{x}_n (f(\mathbf{x}_n) - y_n)$$

**GD update:**

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \lambda \mathbf{X}^T (\mathbf{X} \mathbf{w}^{(t)} - \mathbf{y})$$

For a single weight,

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \lambda \sum_n x_{nj} (f(\mathbf{x}_n) - y_n)$$

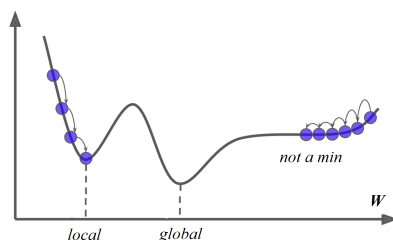
The algorithm uses all training points on each iteration. The algorithm is called *batch gradient descent*.

()

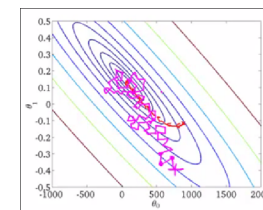
January 29, 2019 12 / 50

### There two main challenges with GD:

- it may converge to a local minimum.
- it may not find a minimum at all. "vanishing gradient".



- GD: move a bit in the negative gradient direction.
- SGD: move a bit in a *noisy* negative gradient direction.
- In SGD, we use one training sample at each iteration.
- Need to randomly shuffle the training examples before calculating it.
- SGD is widely used for larger dataset and can be trained in parallel.



## SGD for Linear Regression

### Algorithm:

initialize  $w^{(0)}$   
for each training sample  $n$ :  
  for each weight  $j$ :

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \lambda x_{nj} (f(\mathbf{x}_n) - y_n)$$

The term stochastic comes from the fact that the gradient is based on a single training sample.

SGD makes progress with each training example as it looks at.

Key point: it could be *much faster to obtain a stochastic gradient!*

## GD versus SGD

In GD we calculate the gradient for all training points

In SGD we calculate the gradient on small batches of training data

In SGD you might not be taking the most optimal route to get to the solution.

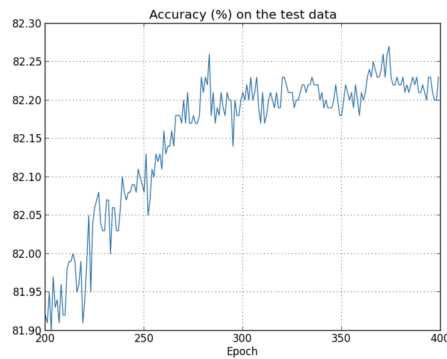
SGD may work for non-convex functions

In SGD you need to go through all batches (the training set) several times (this is called an *epoch*).

You must specify the batch size (a typical size is 256) and number of epochs (a hyperparameter) for a learning algorithm.

## Epoch and overfitting

This shows how test accuracy is changing due to the number of epochs:



()

January 29, 2019 17 / 50

## Outline

- 1 Gradient Descent
- 2 Linear Classifier and Surrogate Losses
- 3 Perceptron
- 4 Logistic Regression

()

January 29, 2019 18 / 50

## General idea to provide ML algorithms

1. Pick a set of **models**  $\mathcal{F}$ 
  - e.g.  $\mathcal{F} = \{f(x) = \mathbf{w}^T \mathbf{x} \mid \mathbf{w} \in \mathbb{R}^D\}$
2. Define **error/loss**  $L(y', y)$
3. Find **empirical risk minimizer (ERM)**:

$$\mathbf{f}^* = \operatorname{argmin}_{\mathbf{f} \in \mathcal{F}} \sum_{n=1}^N L(\mathbf{f}(x_n), y_n)$$

()

January 29, 2019 19 / 50

## Deriving classification algorithms

Let's follow the steps:

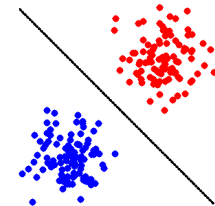
**Step 1.** Pick a set of models  $\mathcal{F}$ .

Again try linear models, but how to predict a label using  $\mathbf{w}^T \mathbf{x}$ ?

*Sign* of  $\mathbf{w}^T \mathbf{x}$  predicts the label:

$$\operatorname{sign}(\mathbf{w}^T \mathbf{x}) = \begin{cases} +1 & \text{if } \mathbf{w}^T \mathbf{x} > 0 \\ -1 & \text{if } \mathbf{w}^T \mathbf{x} \leq 0 \end{cases}$$

(Sometimes use  $\operatorname{sgn}$  for sign too.)



()

January 29, 2019 20 / 50

## The models

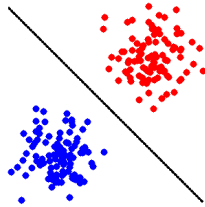
The set of **(separating) hyperplanes**:

$$\mathcal{F} = \{f(\mathbf{x}) = \text{sgn}(\mathbf{w}^T \mathbf{x}) \mid \mathbf{w} \in \mathbb{R}^D\}$$

Good choice for **linearly separable** data, i.e.,  $\exists \mathbf{w}$  s.t.

$$\text{sgn}(\mathbf{w}^T \mathbf{x}_n) = y_n \quad \text{or} \quad y_n \mathbf{w}^T \mathbf{x}_n > 0$$

for all  $n \in [N]$ .

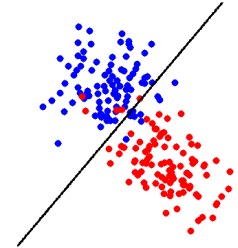


()

January 29, 2019 21 / 50

## The models

Still makes sense for “almost” linearly separable data

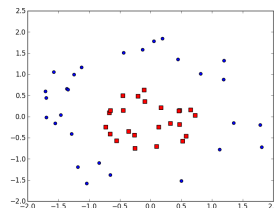
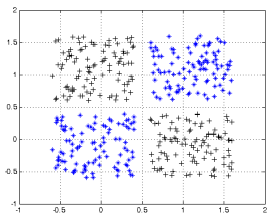


()

January 29, 2019 22 / 50

## The models

For clearly not linearly separable data,



Again can apply a **nonlinear mapping**  $\Phi$ :

$$\mathcal{F} = \{f(\mathbf{x}) = \text{sgn}(\mathbf{w}^T \Phi(\mathbf{x})) \mid \mathbf{w} \in \mathbb{R}^M\}$$

More discussions in the next lectures.

()

January 29, 2019 23 / 50

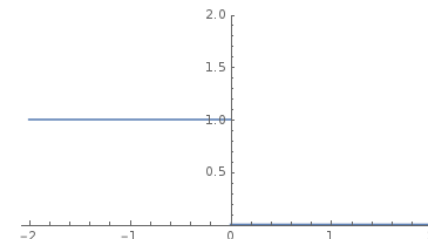
## 0-1 Loss

**Step 2.** Define error/loss  $L(y', y)$ .

Most natural one for classification: **0-1 loss**  $L(y', y) = \mathbb{I}[y' \neq y]$

For classification, more convenient to look at the loss **as a function of**  $y\mathbf{w}^T \mathbf{x}$ . That is, with

$$\ell_{0-1}(z) = \mathbb{I}[z \leq 0]$$



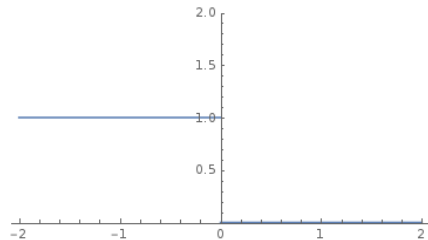
the loss for hyperplane  $\mathbf{w}$  on example  $(\mathbf{x}, y)$  is  $\ell_{0-1}(y\mathbf{w}^T \mathbf{x})$

()

January 29, 2019 24 / 50

## Minimizing 0-1 loss is hard

However, 0-1 loss is *not convex*.



Even worse, minimizing 0-1 loss is *NP-hard in general*.

()

January 29, 2019 25 / 50

## ML becomes convex optimization

**Step 3.** Find ERM:

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^D} \sum_{n=1}^N \ell(y_n \mathbf{w}^T \mathbf{x}_n)$$

where  $\ell(\cdot)$  can be perceptron/hinge/logistic loss

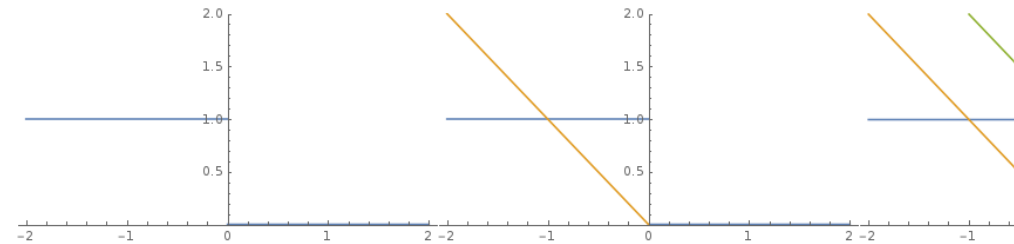
- *no closed-form* in general (unlike linear regression)
- can apply general convex optimization methods

()

January 29, 2019 27 / 50

## Surrogate Losses

Solution: find a **convex surrogate loss**



- **perceptron loss**  $\ell_{\text{perceptron}}(z) = \max\{0, -z\}$  (used in Perceptron)
- **hinge loss**  $\ell_{\text{hinge}}(z) = \max\{0, 1 - z\}$  (used in SVM and many others)
- **logistic loss**  $\ell_{\text{logistic}}(z) = \log(1 + \exp(-z))$  (used in logistic regression)

()

January 29, 2019 26 / 50

## Outline

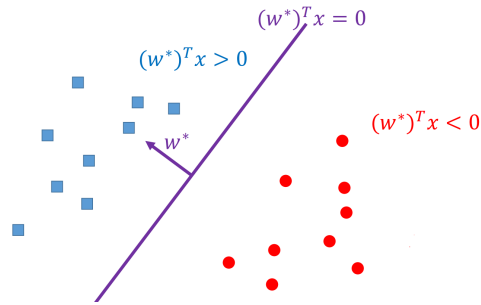
- 1 Gradient Descent
- 2 Linear Classifier and Surrogate Losses
- 3 **Perceptron**
- 4 Logistic Regression

()

January 29, 2019 28 / 50

## The Perceptron

The Perceptron (introduced by Rosenblatt in 1957) is a linear model for classification. Its model is a hyperplane that partitions space into two regions. Perceptron is a rough model for how individual neurons in the brain work.



()

January 29, 2019 29 / 50

## The Perceptron Algorithm

Mathematically: [Stochastic Gradient Descent applied to perceptron loss](#)

i.e. find the minimizer of

$$\begin{aligned} F(\mathbf{w}) &= \sum_{n=1}^N \ell_{\text{perceptron}}(y_n \mathbf{w}^T \mathbf{x}_n) \\ &= \sum_{n=1}^N \max\{0, -y_n \mathbf{w}^T \mathbf{x}_n\} \end{aligned}$$

using SGD

()

January 29, 2019 30 / 50

## Applying GD to perceptron loss

### Objective

$$F(\mathbf{w}) = \sum_{n=1}^N \max\{0, -y_n \mathbf{w}^T \mathbf{x}_n\}$$

Gradient is

$$\nabla F(\mathbf{w}) = \sum_{n=1}^N -\mathbb{I}[y_n \mathbf{w}^T \mathbf{x}_n \leq 0] y_n \mathbf{x}_n$$

(only misclassified examples contribute to the gradient)

### GD update

$$\mathbf{w} \leftarrow \mathbf{w} + \lambda \sum_{n=1}^N \mathbb{I}[y_n \mathbf{w}^T \mathbf{x}_n \leq 0] y_n \mathbf{x}_n$$

()

January 29, 2019 31 / 50

## Applying SGD to perceptron loss

How to construct a stochastic gradient?

### SGD update

$$\mathbf{w} \leftarrow \mathbf{w} + \lambda \mathbb{I}[y_n \mathbf{w}^T \mathbf{x}_n \leq 0] y_n \mathbf{x}_n$$

()

January 29, 2019 32 / 50



## The Perceptron Algorithm

Perceptron algorithm is *SGD with  $\lambda = 1$  applied to perceptron loss*:

Repeat:

- Pick a data point  $x_n$  uniformly at random
- If  $\text{sgn}(\mathbf{w}^T \mathbf{x}_n) \neq y_n$   
 $\mathbf{w} \leftarrow \mathbf{w} + y_n \mathbf{x}_n$

Note:

- The algorithm is online and error driven.
- If the prediction is correct, it does nothing.
- $\mathbf{w}$  is always a *linear combination* of the training examples.
- It uses epochs as a hyperparameter.

()

January 29, 2019 33 / 50

## Why does it make sense?

If the current weight  $\mathbf{w}$  makes a mistake

$$y_n \mathbf{w}^T \mathbf{x}_n < 0$$

then after the update  $\mathbf{w}' = \mathbf{w} + y_n \mathbf{x}_n$  we have

$$y_n \mathbf{w}'^T \mathbf{x}_n = y_n \mathbf{w}^T \mathbf{x}_n + y_n^2 \mathbf{x}_n^T \mathbf{x}_n = y_n \mathbf{w}^T \mathbf{x}_n + \|\mathbf{x}_n\|^2 \geq y_n \mathbf{w}^T \mathbf{x}_n$$

Thus it is more likely to get it right after the update.

()

January 29, 2019 34 / 50

## Any theory?

- If training set is linearly separable, Perceptron *converges in a finite number of steps*
- How long does it take to converge?
- By "how long", what we really mean is "how many updates".
- One way to make this definition is through the notion of margin.
- The margin is the distance between the hyperplane and the nearest point.

()

January 29, 2019 35 / 50

## Perceptron Convergence Theorem

Suppose the perceptron algorithm is run on a linearly separable data set  $\mathcal{D}$  with margin  $\gamma \geq 0$ . Assume that  $\|\mathbf{x}\| = 1$ . Then the algorithm will converge after at most  $\frac{1}{\gamma^2}$  updates.

()

January 29, 2019 36 / 50

## Outline

- 1 Gradient Descent
- 2 Linear Classifier and Surrogate Losses
- 3 Perceptron
- 4 **Logistic Regression**
  - A Probabilistic View
  - Optimization

()

January 29, 2019 37 / 50

## Predicting probability

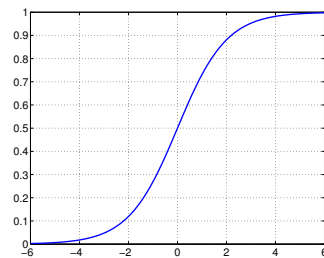
Instead of predicting a discrete label, can we *predict the probability of each label?* i.e. regress the probabilities

One way: **sigmoid function + linear model**

$$\mathbb{P}(y = +1 \mid \mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})$$

where  $\sigma$  is the sigmoid function:

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



()

January 29, 2019 39 / 50

## A simple view

**In one sentence:** find the minimizer of

$$\begin{aligned} F(\mathbf{w}) &= \sum_{n=1}^N \ell_{\text{logistic}}(y_n \mathbf{w}^T \mathbf{x}_n) \\ &= \sum_{n=1}^N \ln(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n}) \end{aligned}$$

*But why logistic loss? and why "regression"?*

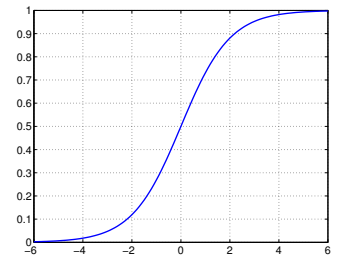
()

January 29, 2019 38 / 50

## Properties

**Properties** of sigmoid  $\sigma(z) = \frac{1}{1+e^{-z}}$

- between 0 and 1 (good as probability)
- $\sigma(\mathbf{w}^T \mathbf{x}) \geq 0.5 \Leftrightarrow \mathbf{w}^T \mathbf{x} \geq 0$ , consistent with predicting the label with  $\text{sgn}(\mathbf{w}^T \mathbf{x})$
- larger  $\mathbf{w}^T \mathbf{x} \Rightarrow$  larger  $\sigma(\mathbf{w}^T \mathbf{x}) \Rightarrow$  higher *confidence* in label 1
- $\sigma(z) + \sigma(-z) = 1$  for all  $z$



The probability of label  $-1$  is naturally

$$1 - \mathbb{P}(y = +1 \mid \mathbf{x}; \mathbf{w}) = 1 - \sigma(\mathbf{w}^T \mathbf{x}) = \sigma(-\mathbf{w}^T \mathbf{x})$$

and thus

$$\mathbb{P}(y \mid \mathbf{x}; \mathbf{w}) = \sigma(y \mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-y \mathbf{w}^T \mathbf{x}}}$$

()

January 29, 2019 40 / 50

## How to regress with discrete labels?

*What we observe are labels, not probabilities.*

Take a **probabilistic view**

- assume data is generated in this way by some  $\mathbf{w}$
- perform Maximum Likelihood Estimation (MLE)

Specifically, what is the probability of seeing label  $y_1, \dots, y_n$  given  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , as a function of some  $\mathbf{w}$ ?

$$P(\mathbf{w}) = \prod_{n=1}^N \mathbb{P}(y_n | \mathbf{x}_n; \mathbf{w})$$

**MLE**: find  $\mathbf{w}^*$  that **maximizes the probability**  $P(\mathbf{w})$

## The MLE solution

$$\begin{aligned}
 \mathbf{w}^* &= \operatorname{argmax}_{\mathbf{w}} P(\mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} \prod_{n=1}^N \mathbb{P}(y_n | \mathbf{x}_n; \mathbf{w}) \\
 &= \operatorname{argmax}_{\mathbf{w}} \sum_{n=1}^N \ln \mathbb{P}(y_n | \mathbf{x}_n; \mathbf{w}) = \operatorname{argmin}_{\mathbf{w}} \sum_{n=1}^N -\ln \mathbb{P}(y_n | \mathbf{x}_n; \mathbf{w}) \\
 &= \operatorname{argmin}_{\mathbf{w}} \sum_{n=1}^N \ln(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n}) = \operatorname{argmin}_{\mathbf{w}} \sum_{n=1}^N \ell_{\text{logistic}}(y_n \mathbf{w}^T \mathbf{x}_n) \\
 &= \operatorname{argmin}_{\mathbf{w}} F(\mathbf{w})
 \end{aligned}$$

i.e. *minimizing logistic loss is exactly doing MLE for the sigmoid model!*

()

January 29, 2019 41 / 50

()

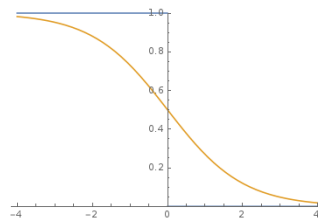
January 29, 2019 42 / 50

## Let's apply SGD again

$$\begin{aligned}
 \mathbf{w} &\leftarrow \mathbf{w} - \lambda \nabla F(\mathbf{w}) \\
 &= \mathbf{w} - \lambda \nabla_{\mathbf{w}} \ell_{\text{logistic}}(y_n \mathbf{w}^T \mathbf{x}_n) \\
 &= \mathbf{w} - \lambda \left( \frac{\partial \ell_{\text{logistic}}(z)}{\partial z} \Big|_{z=y_n \mathbf{w}^T \mathbf{x}_n} \right) y_n \mathbf{x}_n \\
 &= \mathbf{w} - \lambda \left( \frac{-e^{-z}}{1 + e^{-z}} \Big|_{z=y_n \mathbf{w}^T \mathbf{x}_n} \right) y_n \mathbf{x}_n \\
 &= \mathbf{w} + \lambda \sigma(-y_n \mathbf{w}^T \mathbf{x}_n) y_n \mathbf{x}_n \\
 &= \mathbf{w} + \lambda \mathbb{P}(-y_n | \mathbf{x}_n; \mathbf{w}) y_n \mathbf{x}_n
 \end{aligned}$$

This is a *soft version of Perceptron!*

$\mathbb{P}(-y_n | \mathbf{x}_n; \mathbf{w})$  versus  $\mathbb{I}[y_n \neq \operatorname{sgn}(\mathbf{w}^T \mathbf{x}_n)]$



()

January 29, 2019 43 / 50

## Newton method

Newtons method is an extension of steepest descent, where the second-order term in the Taylor series is used.

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2} f''(x_0)(x - x_0)^2$$

Let us minimize the right hand side:

$$f'(x_0) + f''(x_0)(x - x_0) = 0 \quad \text{or} \quad x = x_0 - \frac{f'(x_0)}{f''(x_0)}$$

We will iterate this procedure

$$x_{n+1} = x_n - \frac{f'(x_n)}{f''(x_n)}$$

()

January 29, 2019 44 / 50

## Deriving Newton method

This could be generalized for functions  $f$  of several variables:

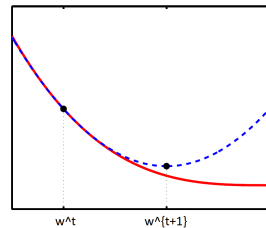
$$x_{n+1} = x_n - \mathbf{H}^{-1}(x_n) \nabla f(x_n)$$

where  $\mathbf{H}$  is the Hessian

$$H_{ij} = \frac{\partial^2 F(\mathbf{x})}{\partial x_i \partial x_j}$$

Therefore, for convex  $F$  (so  $H_t$  is **positive semidefinite**) we obtain **Newton method**:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \mathbf{H}_t^{-1} \nabla F(\mathbf{w}^{(t)})$$



()

January 29, 2019 45 / 50

## Comparing GD and Newton

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \lambda \nabla F(\mathbf{w}^{(t)}) \quad (\text{GD})$$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \mathbf{H}_t^{-1} \nabla F(\mathbf{w}^{(t)}) \quad (\text{Newton})$$

Both are iterative optimization procedures, but Newton method

- has no learning rate  $\lambda$  (**so no tuning needed!**)
- converges **super fast** in terms of #iterations needed
- requires **second-order** information

()

January 29, 2019 46 / 50

## Applying Newton to logistic loss

$$\nabla_{\mathbf{w}} \ell_{\text{logistic}}(y_n \mathbf{w}^T \mathbf{x}_n) = -\sigma(-y_n \mathbf{w}^T \mathbf{x}_n) y_n \mathbf{x}_n$$

$$\begin{aligned} \nabla_{\mathbf{w}}^2 \ell_{\text{logistic}}(y_n \mathbf{w}^T \mathbf{x}_n) &= \left( \frac{\partial \sigma(z)}{\partial z} \Big|_{z=-y_n \mathbf{w}^T \mathbf{x}_n} \right) y_n^2 \mathbf{x}_n \mathbf{x}_n^T \\ &= \left( \frac{e^{-z}}{(1 + e^{-z})^2} \Big|_{z=-y_n \mathbf{w}^T \mathbf{x}_n} \right) \mathbf{x}_n \mathbf{x}_n^T \\ &= \sigma(y_n \mathbf{w}^T \mathbf{x}_n) (1 - \sigma(y_n \mathbf{w}^T \mathbf{x}_n)) \mathbf{x}_n \mathbf{x}_n^T \end{aligned}$$

### Exercises:

- why is the Hessian of logistic loss positive semidefinite?
- can we apply Newton method to perceptron/hinge loss?

()

January 29, 2019 47 / 50

## Summary

Linear models for classification:

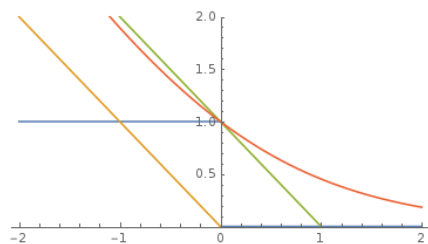
Step 1. Model is the set of **separating hyperplanes**

$$\mathcal{F} = \{f(\mathbf{x}) = \text{sgn}(\mathbf{w}^T \mathbf{x}) \mid \mathbf{w} \in \mathbb{R}^D\}$$

()

January 29, 2019 48 / 50

Step 2. Pick the **surrogate loss**



- **perceptron loss**  $\ell_{\text{perceptron}}(z) = \max\{0, -z\}$  (used in Perceptron)
- **hinge loss**  $\ell_{\text{hinge}}(z) = \max\{0, 1 - z\}$  (used in SVM and many others)
- **logistic loss**  $\ell_{\text{logistic}}(z) = \log(1 + \exp(-z))$  (used in logistic regression)

Step 3. Find empirical risk minimizer (ERM):

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^D} \sum_{n=1}^N \ell(y_n \mathbf{w}^T \mathbf{x}_n)$$

using **GD/SGD/Newton**.