# Machine Learning (CSE 446):
# Learning as Minimizing Loss (continued)

Noah Smith
© 2017

University of Washington
nasmith@cs.washington.edu

October 25, 2017

warning

## Gradient Descent

**Data**: function $F : \mathbb{R}^d \to \mathbb{R}$, number of iterations $K$, step sizes $\langle \eta^{(1)}, \ldots, \eta^{(K)} \rangle$
**Result**: $\mathbf{z} \in \mathbb{R}^d$
initialize: $\mathbf{z}^{(0)} = \mathbf{0}$;
**for** $k \in \{1, \ldots, K\}$ **do**
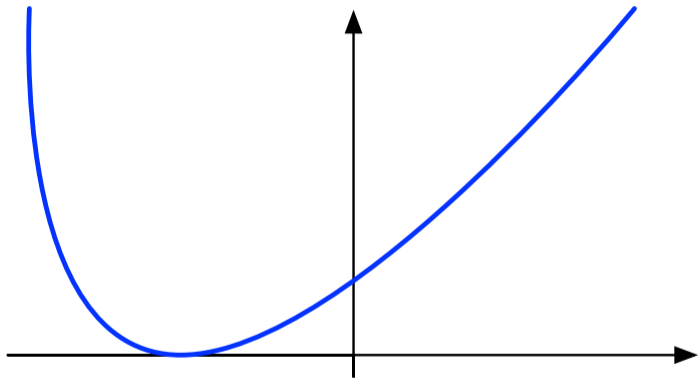$\quad$ $\mathbf{g}^{(k)} = \nabla_{\mathbf{z}} F(\mathbf{z}^{(k-1)})$;
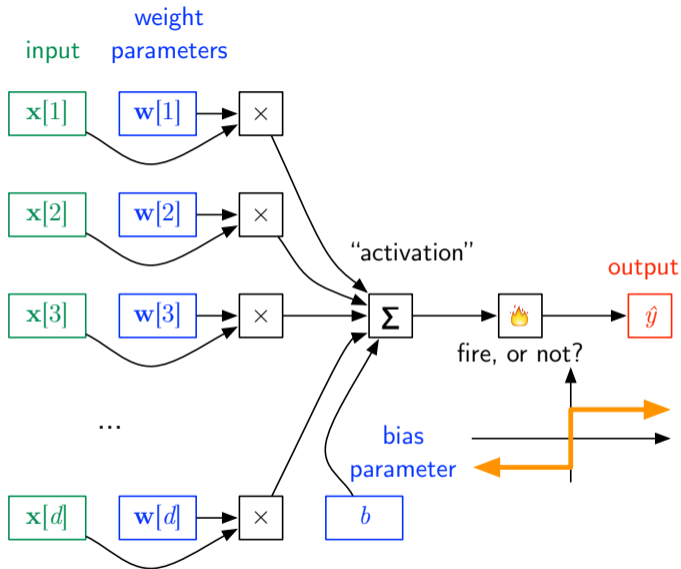$\quad$ $\mathbf{z}^{(k)} = \mathbf{z}^{(k-1)} - \eta^{(k)} \cdot \mathbf{g}^{(k)}$;
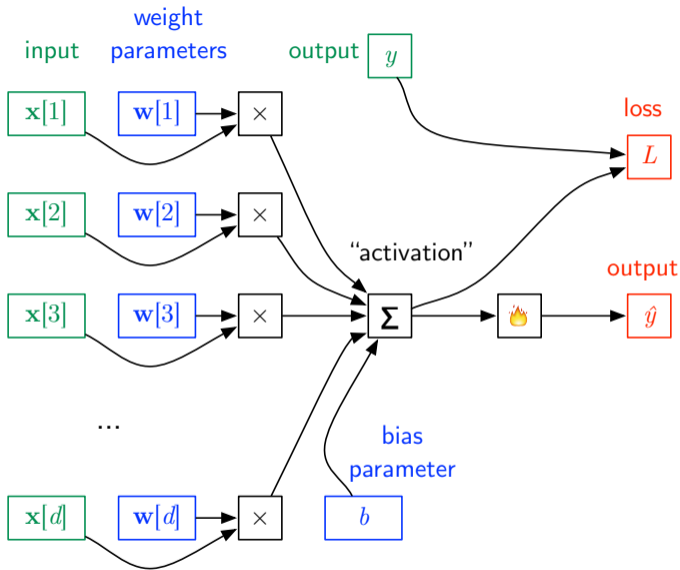**end**
return $\mathbf{z}^{(K)}$;

**Algorithm 1:** GRADIENTDESCENT

# Gradient Descent

weight parameters

input

$\mathbf{x}[1]$ $\mathbf{w}[1]$ $\times$

$\mathbf{x}[2]$ $\mathbf{w}[2]$ $\times$

$\mathbf{x}[3]$ $\mathbf{w}[3]$ $\times$

"activation"

$\Sigma$ 🔥 $\hat{y}$

fire, or not?

output

...

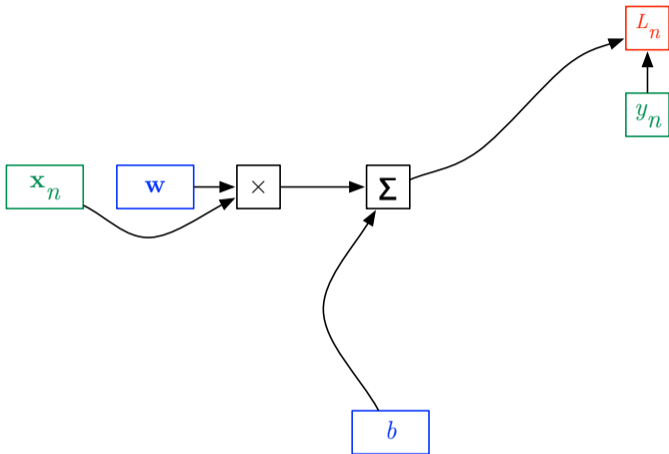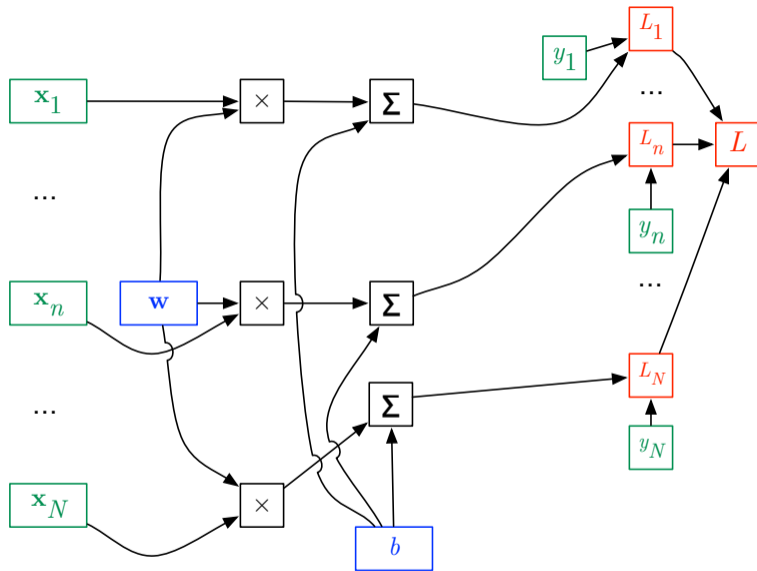$\mathbf{x}[d]$ $\mathbf{w}[d]$ $\times$

bias parameter

$b$

Classification
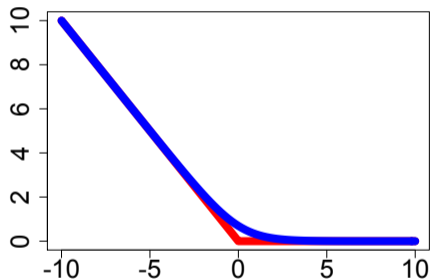
Training (one example)

Training (one example)

Training ($D$)

## Log Loss

The log loss is continuous, convex, and differentiable. It is closely related to the perceptron loss.



$$L_n(\mathbf{w}) = \log\left(1 + \exp\left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)\right)$$

# First Derivative of the Log Loss for One Example

$$\frac{\partial L_n}{\partial w[j]} = \frac{\partial}{\partial w[j]} \log \left( 1 + \exp \left( -y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b) \right) \right)$$

# First Derivative of the Log Loss for One Example

$$\frac{\partial L_n}{\partial w[j]} = \frac{\partial}{\partial w[j]} \log \left(1 + \exp\left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)\right)$$

$$= \frac{1}{1 + \exp\left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot \frac{\partial}{\partial w[j]} \left[1 + \exp\left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)\right]$$

Remember: $\frac{\partial}{\partial x} \log f(x) = \frac{\frac{\partial}{\partial x} f(x)}{f(x)}$.

# First Derivative of the Log Loss for One Example

$$\frac{\partial L_n}{\partial w[j]} = \frac{\partial}{\partial w[j]} \log \left(1 + \exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)\right)$$

$$= \frac{1}{1 + \exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot \frac{\partial}{\partial w[j]} \left[1 + \exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)\right]$$

$$= \frac{1}{1 + \exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot \frac{\partial}{\partial w[j]} \exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)$$

# First Derivative of the Log Loss for One Example

$$\frac{\partial L_n}{\partial w[j]} = \frac{\partial}{\partial w[j]} \log\left(1 + \exp\left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)\right)$$

$$= \frac{1}{1 + \exp\left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot \frac{\partial}{\partial w[j]} \left[1 + \exp\left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)\right]$$

$$= \frac{1}{1 + \exp\left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot \frac{\partial}{\partial w[j]} \exp\left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)$$

$$= \frac{\exp\left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)}{1 + \exp\left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot \frac{\partial}{\partial w[j]} \left[-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right]$$

Remember: $\frac{\partial}{\partial x} \exp f(x) = \exp f(x) \cdot \frac{\partial}{\partial x} f(x)$.

# First Derivative of the Log Loss for One Example

$$\frac{\partial L_n}{\partial w[j]} = \frac{\partial}{\partial w[j]} \log \left(1 + \exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)\right)$$

$$= \frac{1}{1 + \exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot \frac{\partial}{\partial w[j]} \left[1 + \exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)\right]$$

$$= \frac{1}{1 + \exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot \frac{\partial}{\partial w[j]} \exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)$$

$$= \frac{\exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)}{1 + \exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot \frac{\partial}{\partial w[j]} \left[-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right]$$

$$= \frac{\exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)}{1 + \exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot -y_n \cdot \mathbf{x}_n[j]$$

# First Derivative of the Log Loss for One Example

$$\frac{\partial L_n}{\partial w[j]} = \frac{\partial}{\partial w[j]} \log \left(1 + \exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)\right)$$

$$= \frac{1}{1 + \exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot \frac{\partial}{\partial w[j]} \left[1 + \exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)\right]$$

$$= \frac{1}{1 + \exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot \frac{\partial}{\partial w[j]} \exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)$$

$$= \frac{\exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)}{1 + \exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot \frac{\partial}{\partial w[j]} \left[-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right]$$

$$= \frac{\exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)}{1 + \exp \left(-y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot -y_n \cdot \mathbf{x}_n[j]$$

$$= \frac{1}{1 + \exp \left(y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot -y_n \cdot \mathbf{x}_n[j]$$

$\frac{\exp x}{1 + \exp x} = \frac{\exp x}{1 + \exp x} \cdot \frac{\exp -x}{\exp -x} = \frac{1}{1 + \exp -x}.$

## Gradient of the Log Loss, All Examples

One example:

$$\frac{1}{1 + \exp\left(y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot -y_n \cdot \mathbf{x}_n[j]$$

Sum over all examples:

$$\sum_{n=1}^{N} \frac{1}{1 + \exp\left(y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot -y_n \cdot \mathbf{x}_n[j]$$

Gradient stacks all first derivatives into a vector:

$$\sum_{n=1}^{N} \frac{1}{1 + \exp\left(y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot -y_n \cdot \mathbf{x}_n$$

# Perceptron Learning vs. Gradient Descent on Log Loss

Updating rule for the perceptron (single example $n$):

$$\mathbf{w} \leftarrow \mathbf{w} + [\![\hat{y}_n \neq y_n]\!] \cdot y_n \cdot \mathbf{x}_n$$
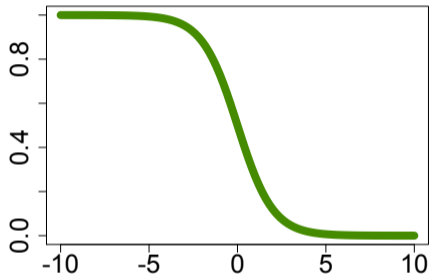
# Perceptron Learning vs. Gradient Descent on Log Loss

Updating rule for the perceptron (single example $n$):

$$\mathbf{w} \leftarrow \mathbf{w} + [\![\hat{y}_n \neq y_n]\!] \cdot y_n \cdot \mathbf{x}_n$$
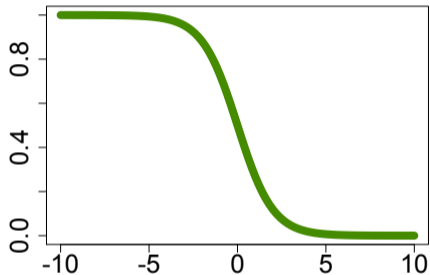
Gradient descent on log loss:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \cdot \sum_{n=1}^{N} \frac{1}{1 + \exp\left(y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot -y_n \cdot \mathbf{x}_n$$

$$= \mathbf{w} + \eta \cdot \sum_{n=1}^{N} \frac{1}{1 + \exp\left(y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot y_n \cdot \mathbf{x}_n$$

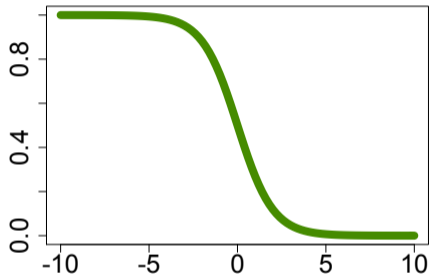# An Interesting Function



$$\sigma(x) = \frac{1}{1 + \exp x}$$

# An Interesting Function



What we plug in to this function is
$y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)$.

$$\sigma(x) = \frac{1}{1 + \exp x}$$

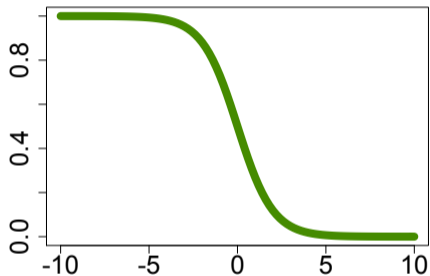# An Interesting Function



What we plug in to this function is
$y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)$.
If example $n$ is getting classified correctly,
$\sigma(y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b))$ goes to zero.

$$\sigma(x) = \frac{1}{1 + \exp x}$$

# An Interesting Function



What we plug in to this function is
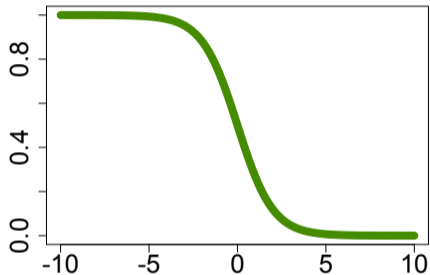$y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)$.
If example $n$ is getting classified correctly,
$\sigma(y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b))$ goes to zero.
If example $n$ is getting classified incorrectly,
$\sigma(y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b))$ goes to one.

$$\sigma(x) = \frac{1}{1 + \exp x}$$

# An Interesting Function



$$\sigma(x) = \frac{1}{1 + \exp x}$$

What we plug in to this function is
$y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)$.
If example $n$ is getting classified correctly,
$\sigma(y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b))$ goes to zero.
If example $n$ is getting classified incorrectly,
$\sigma(y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b))$ goes to one.
Not "is it wrong?" but instead "how wrong is it?"

# Perceptron Learning vs. Gradient Descent on Log Loss

Updating rule for the perceptron (single example $n$):

$$\mathbf{w} \leftarrow \mathbf{w} + [\![\hat{y}_n \neq y_n]\!] \cdot y_n \cdot \mathbf{x}_n$$

Gradient descent on log loss:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \cdot \sum_{n=1}^{N} \frac{1}{1 + \exp\left(y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot -y_n \cdot \mathbf{x}_n$$

$$= \mathbf{w} + \eta \cdot \sum_{n=1}^{N} \frac{1}{1 + \exp\left(y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)} \cdot y_n \cdot \mathbf{x}_n$$

$$= \mathbf{w} + \sum_{n=1}^{N} \underbrace{\frac{\eta}{1 + \exp\left(y_n \cdot (\mathbf{w} \cdot \mathbf{x}_n + b)\right)}}_{\text{``soft'' error test}} \cdot y_n \cdot \mathbf{x}_n$$

# Logistic Regression

Classifier:

$$f(\mathbf{x}) = \text{sign} \left( \mathbf{w} \cdot \mathbf{x} + b \right)$$

(Identical to the neuron-inspired classifier we trained using the perceptron!)

To learn, solve:

$$\min_{\mathbf{w}, b} \frac{1}{N} \sum_{n=1}^{N} \underbrace{\log \left( 1 + \exp \left( -y_n \cdot \left( \mathbf{w} \cdot \mathbf{x}_n + b \right) \right) \right)}_{\text{log loss}} + R(\mathbf{w}, b)$$

# Linear Regression

For *continuous* output $y$.

Predictor:

$$f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$$

To learn, solve:

$$\min_{\mathbf{w}, b} \frac{1}{N} \sum_{n=1}^{N} \underbrace{(y_n - (\mathbf{w} \cdot \mathbf{x}_n + b))^2}_{\text{squared loss}} + R(\mathbf{w}, b)$$

# Coming Soon

- Gradient descent is usually too slow; we can do better.
- Is there a deeper connection between perceptron learning and gradient descent?
- More loss functions!
- More activation functions!
- More regularization functions!
- A probabilistic interpretation of logistic regression and linear regression.