

# Logistic Regression and Convex Optimization

Anthony Platanios

# Logistic Regression

Binary:

$$P(Y = 1 \mid X = \mathbf{x}) = \frac{1}{1 + e^{\mathbf{w}^\top \mathbf{x}}}$$

Multiclass:

$$P(Y = y_k \mid X = \mathbf{x}) = \frac{e^{\mathbf{w}_k^\top \mathbf{x}}}{1 + \sum_{k'=1}^{K-1} e^{\mathbf{w}_{k'}^\top \mathbf{x}}}$$

# Logistic Regression Decision Boundary

Linear classification boundary, but why?

$$P(Y = 1 | X = \mathbf{x}) = P(Y = 0 | X = \mathbf{x}) \Rightarrow$$

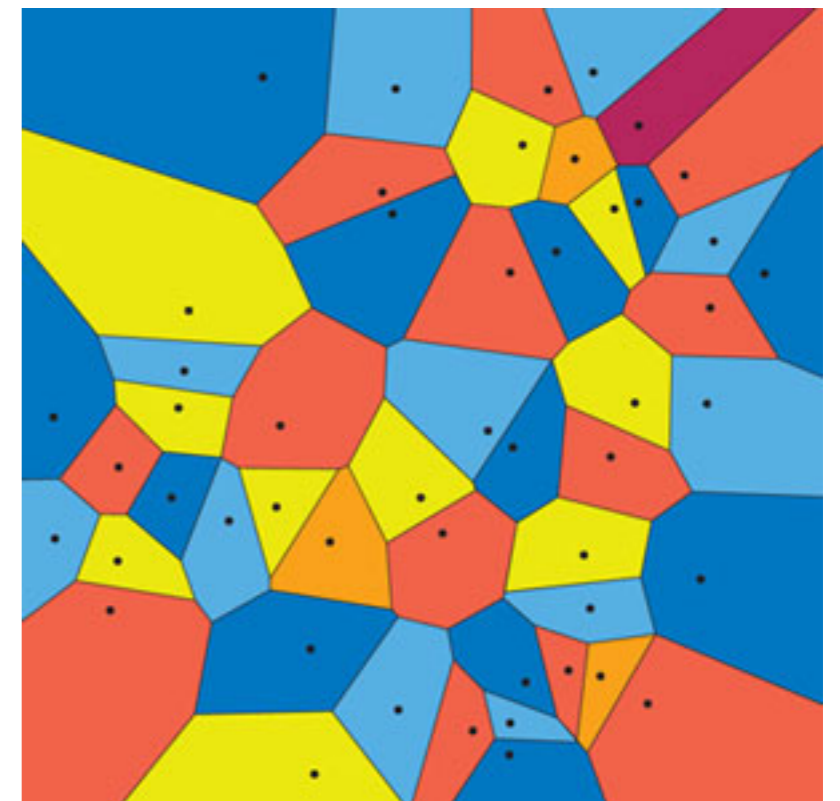
$$\frac{1}{1 + e^{w^\top \mathbf{x}}} = \frac{e^{w^\top \mathbf{x}}}{1 + e^{w^\top \mathbf{x}}} \Rightarrow$$

$$w^\top \mathbf{x} = 0$$



**Hyperplane**

Multiclass: **Voronoi**



# Training Logistic Regression

We use maximum condition likelihood estimation (MCLE):

$$\mathbf{w} = \arg \max_{\mathbf{w}} \prod_{i=1}^n P(Y = y^i \mid X = \mathbf{x}^i, \mathbf{w}) \Rightarrow$$

$$\mathbf{w} = \arg \max_{\mathbf{w}} \sum_{i=1}^n y^i (\mathbf{w}^\top \mathbf{x}^i) - \log (1 + e^{\mathbf{w}^\top \mathbf{x}^i})$$

But how do we solve this?

**Convex Optimization**

# Gradient Descent

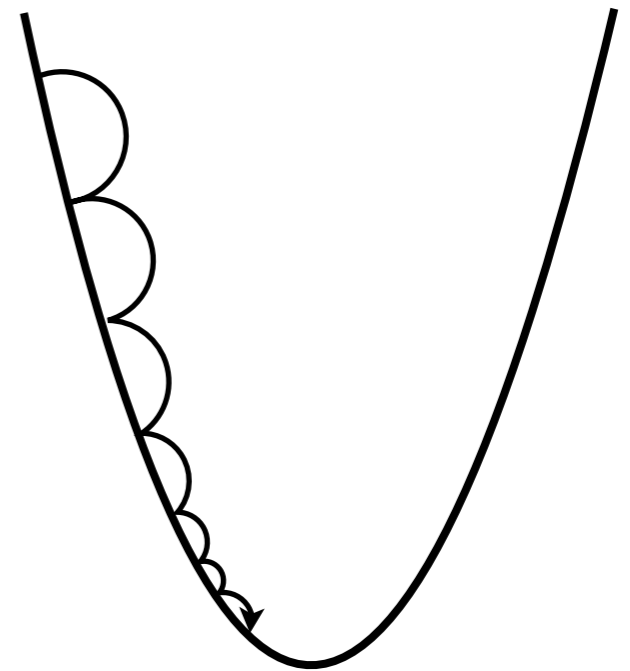
Problem:

$$\arg \min_w f(w)$$

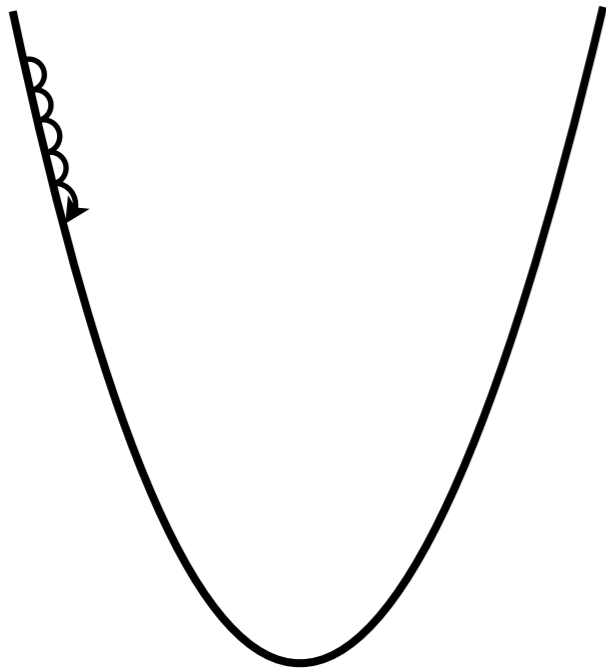
Iterative Solution:

$$w \leftarrow w - \eta \frac{\partial f(w)}{\partial w}$$

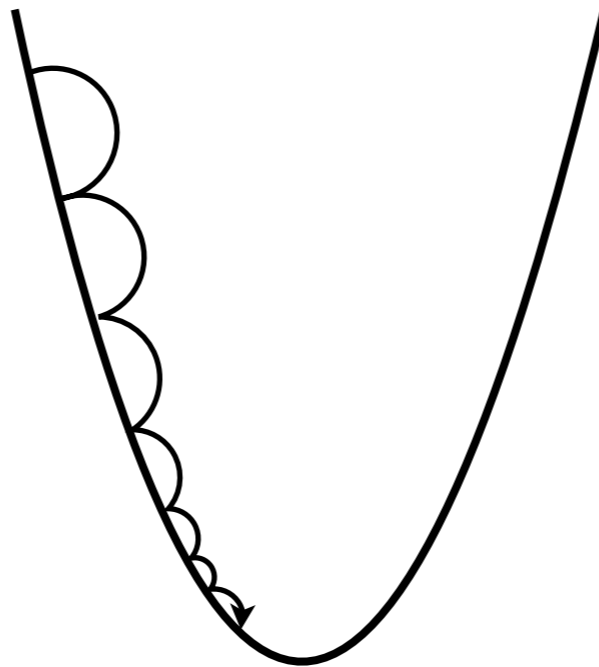
↑  
**Step Size**



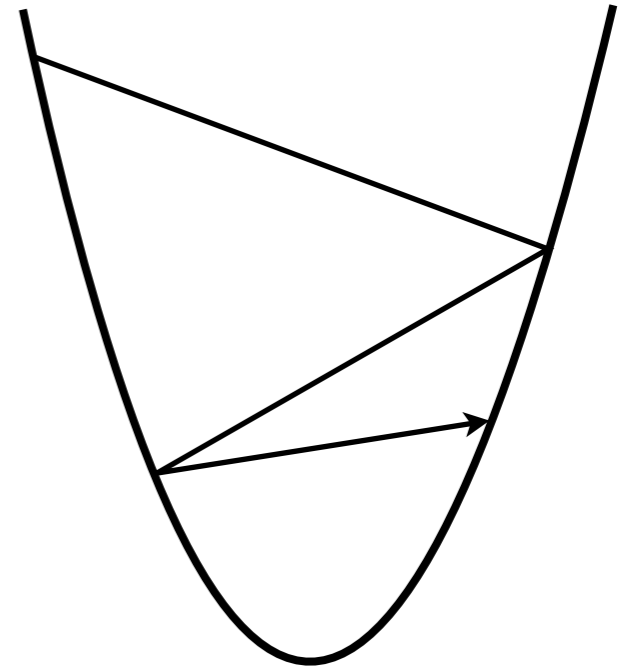
# Step Size Selection



Too Small



**Just Right**



Too Big

# Step Size Selection

Exact Line Search:  $\eta^* = \arg \min_{\eta \geq 0} f \left( \mathbf{w} - \eta \frac{\partial f(\mathbf{w})}{\partial \mathbf{w}} \right)$

Most often we cannot obtain a closed-form solution!

Backtracking Line Search: Choose  $0 < \beta < 1$ ,  $0 < \alpha < 0.5$  and  $\eta_0 \geq 0$  and while:

$$f \left( \mathbf{w} - \eta_t \frac{\partial f(\mathbf{w})}{\partial \mathbf{w}} \right) > f(\mathbf{w}) - \alpha \eta_t \left\| \frac{\partial f(\mathbf{w})}{\partial \mathbf{w}} \right\|_2^2$$

set:  $\eta_{t+1} = \beta \eta_t$

We start with a big step size and keep decreasing it until the update generates a sufficient decrease for our function. We want our next iterate to beat the criterion value of a linear approximation of our function at the current point. The above inequality is called the **Armijo rule** and it is often used in combination with a **curvature condition**. Check the [wikipedia page on Wolfe conditions](#) for more information.

# Other Convex Optimization Methods

There are other convex optimization methods that use even second derivative information, such as **Newton's method**.



Use the **inverse Hessian** as the step size  
(makes use of **curvature** information)

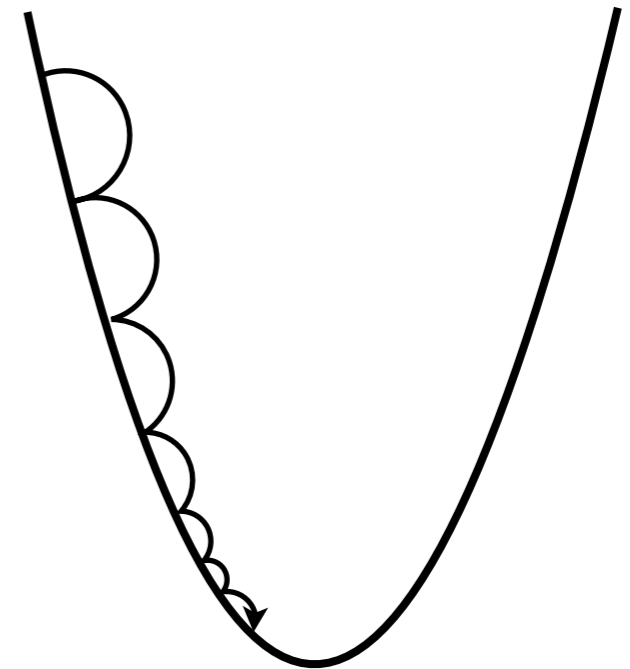
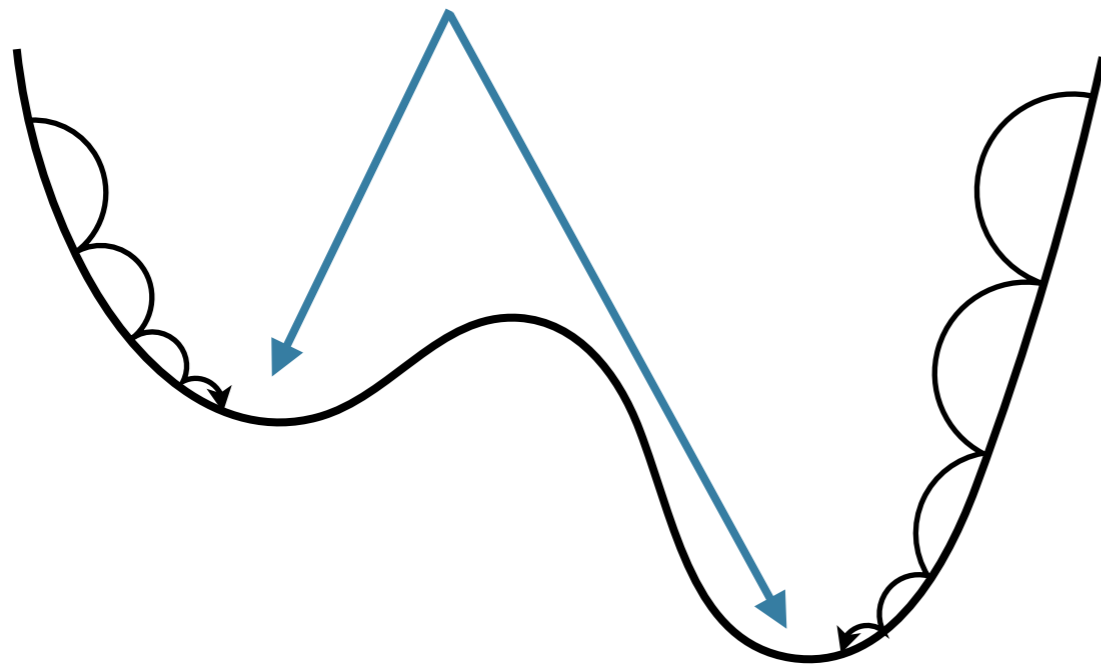
If the Hessian is too expensive to compute, people might also use **quasi-Newton** methods such as the well known **LBFGS** algorithm.

But, what is **convexity**?



# Convexity

**Local Minima**



**Non-convex**

**Convex**

$\forall \mathbf{w}, \mathbf{u} \in \text{dom}(f)$  and  $\forall \theta \in [0, 1]$  :

$$f(\theta \mathbf{w} + (1 - \theta) \mathbf{u}) \leq \theta f(\mathbf{w}) + (1 - \theta) f(\mathbf{u})$$

# Operations that Preserve Convexity

Nonnegative linear combinations: If  $f_1, \dots, f_m$  are convex, then  $\alpha_1 f_1 + \dots + \alpha_m f_m$  is convex for any  $\alpha_1, \dots, \alpha_m \geq 0$ .

Affine compositions: If  $f$  is convex, then  $g(\mathbf{w}) = f(\mathbf{A}\mathbf{w} + \mathbf{b})$  is also convex.

Pointwise maximum: If  $f_1, \dots, f_m$  are convex, then  $\max \{f_1, \dots, f_m\}$  is also convex.

# Back to Logistic Regression

Sum of **affine functions** and **convex functions**.


$$\boldsymbol{w} = \arg \max_{\boldsymbol{w}} L(\boldsymbol{w}) = \arg \max_{\boldsymbol{w}} \sum_{i=1}^n y^i (\boldsymbol{w}^\top \boldsymbol{x}^i) - \log (1 + e^{\boldsymbol{w}^\top \boldsymbol{x}^i})$$

$$\frac{\partial L(\boldsymbol{w})}{\partial \boldsymbol{w}} = \sum_{i=1}^n \boldsymbol{x}^i \left( y^i - \frac{e^{\boldsymbol{w}^\top \boldsymbol{x}^i}}{1 + e^{\boldsymbol{w}^\top \boldsymbol{x}^i}} \right)$$

$$\boldsymbol{w} \leftarrow \boldsymbol{w} + \eta \frac{\partial L(\boldsymbol{w})}{\partial \boldsymbol{w}}$$

# Logistic Regression with a Prior

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

$$\|\mathbf{w}\|_2^2 \triangleq \sum_j w_j^2$$


$$L(\mathbf{w}) = \sum_{i=1}^n y^i (\mathbf{w}^\top \mathbf{x}^i) - \log(1 + e^{\mathbf{w}^\top \mathbf{x}^i}) - \frac{1}{2\sigma^2} \|\mathbf{w}\|_2^2$$

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \sum_{i=1}^n \mathbf{x}^i \left( y^i - \frac{e^{\mathbf{w}^\top \mathbf{x}^i}}{1 + e^{\mathbf{w}^\top \mathbf{x}^i}} \right) - \frac{1}{\sigma^2} \mathbf{w}$$

More generally we call this **L2 regularization**:

$$\arg \min_{\mathbf{w}} f(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$$

# Regularization in Optimization

More generally we have **Lp regularization**:

$$\arg \min_{\mathbf{w}} f(\mathbf{w}) + \lambda \|\mathbf{w}\|_p^p$$

**L1** is a special case, frequently used in practice to induce **sparsity** in the solution:

$$\arg \min_{\mathbf{w}} f(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$$

# Stochastic Gradient Descent

We can randomly sample terms of that sum and get an estimate of the gradient in order to speed things up



$$\boldsymbol{w} = \arg \max_{\boldsymbol{w}} L(\boldsymbol{w}) = \arg \max_{\boldsymbol{w}} \sum_{i=1}^n y^i (\boldsymbol{w}^\top \boldsymbol{x}^i) - \log (1 + e^{\boldsymbol{w}^\top \boldsymbol{x}^i})$$

$$\frac{\partial L(\boldsymbol{w})}{\partial \boldsymbol{w}} = \sum_{i=1}^n \boldsymbol{x}^i \left( y^i - \frac{e^{\boldsymbol{w}^\top \boldsymbol{x}^i}}{1 + e^{\boldsymbol{w}^\top \boldsymbol{x}^i}} \right)$$

$$\boldsymbol{w} \leftarrow \boldsymbol{w} + \eta \frac{\partial L(\boldsymbol{w})}{\partial \boldsymbol{w}}$$