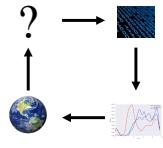
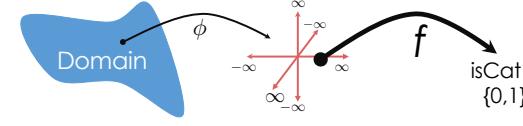


Finishing  
**Logistic Regression**  
Classification with Linear Models

Slides by:  
**Joseph E. Gonzalez**  
[jegonzal@cs.berkeley.edu](mailto:jegonzal@cs.berkeley.edu)

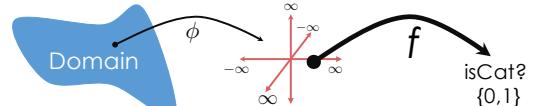


**Classification**



Predicting categorical responses variables

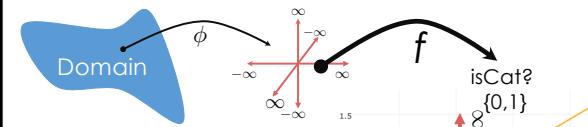
**Classification**



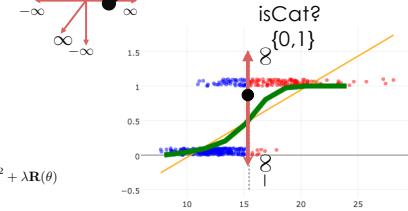
Can we just use least squares?

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n (y_i - f_{\theta}(x_i))^2 + \lambda R(\theta)$$

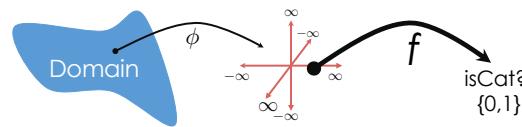
**Classification**



Can we just use least squares?

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n (y_i - f_{\theta}(x_i))^2 + \lambda R(\theta)$$


**Classification**



Can we just use least squares?

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n (y_i - f_{\theta}(x_i))^2 + \lambda R(\theta)$$

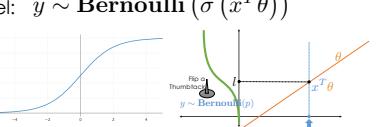
- Yes ... (can be easy to compute)
- Need a decision function (e.g.,  $f(x) > 0.5$ ) ...
- Difficult to interpret model ...

Eh...



**Recap: Logistic Regression**

- Defined a model:  $y \sim \text{Bernoulli}(\sigma(x^T \theta))$

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


- Constructed a likelihood function:

$$\mathcal{L}(\theta) = \prod_{i=1}^n \sigma(x_i^T \theta)^{y_i} (1 - \sigma(x_i^T \theta))^{(1-y_i)}$$

- Constructed a likelihood function:

$$\mathcal{L}(\theta) = \prod_{i=1}^n \sigma(x_i^T \theta)^{y_i} (1 - \sigma(x_i^T \theta))^{(1-y_i)}$$

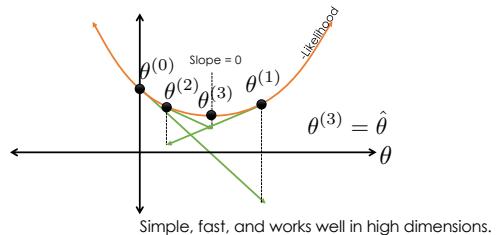
- Computed the gradient of the log-likelihood

$$\nabla_{\theta} \log \mathcal{L}(\theta) = \sum_{i=1}^n (y_i - \sigma(x_i^T \theta)) x_i$$

➤ Set equal to 0 and solve ... no analytic solution ☺

- Gradient Descent

## Gradient Descent Intuition



- Gradient Descent

$$\theta^{(0)} \leftarrow \text{random initial vector}$$

For  $\tau$  from 0 to convergence

$$\theta^{(\tau+1)} \leftarrow \theta^{(\tau)} - \rho(\tau) \nabla_{\theta} (-\log \mathcal{L}(\theta)) \quad \begin{matrix} \text{Evaluated} \\ \text{at} \\ \theta^{(\tau)} \end{matrix}$$

$$\nabla_{\theta} (-\log \mathcal{L}(\theta)) = \sum_{i=1}^n (\sigma(x_i^T \theta) - y_i) x_i$$

Maximum Likelihood → Minimize Negative Log-likelihood

- Gradient Descent

$$\theta^{(0)} \leftarrow \text{random initial vector}$$

For  $\tau$  from 0 to convergence

$$\theta^{(\tau+1)} \leftarrow \theta^{(\tau)} - \rho(\tau) \nabla_{\theta} (-\log \mathcal{L}(\theta)) \quad \begin{matrix} \text{Evaluated} \\ \text{at} \\ \theta^{(\tau)} \end{matrix}$$

- Stochastic Gradient Descent (SGD)

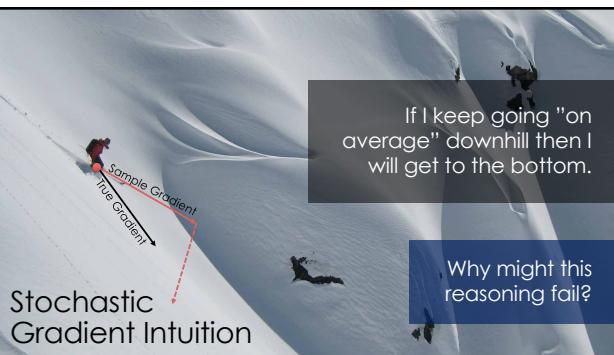
$$\theta^{(0)} \leftarrow \text{random initial vector}$$

For  $\tau$  from 0 to convergence

$$\mathcal{B} \leftarrow \text{select } k \text{ random indices}$$

$$\theta^{(\tau+1)} \leftarrow \theta^{(\tau)} - \rho(\tau) \frac{n}{k} \sum_{i \in \mathcal{B}} \nabla_{\theta} (-\log \mathcal{L}(\theta | x_i, y_i)) \quad \begin{matrix} \text{Eval.} \\ \text{at} \\ \theta^{(\tau)} \end{matrix}$$

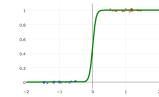
Estimate the gradient by sampling the data



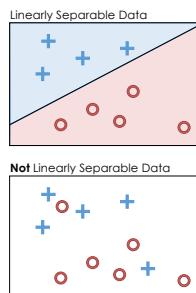
Python Demo!

## Linearly Separable Data

- A classification dataset is said to be linearly separable if there exists a hyperplane that separates the two classes.
- If data is linearly separable, logistic regression requires regularization



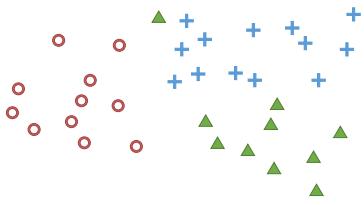
Weights go to infinity!



## Python Demo!

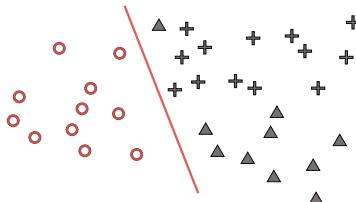
## Multiclass (more than 2) Classification

- **One-vs-rest** train separate binary classifiers for each class



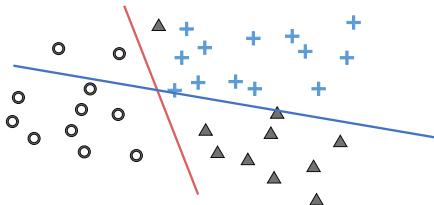
## Multiclass (more than 2) Classification

- **One-vs-rest** train separate binary classifiers for each class



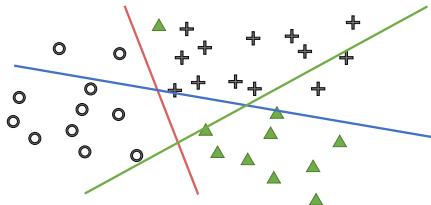
## Multiclass (more than 2) Classification

- **One-vs-rest** train separate binary classifiers for each class



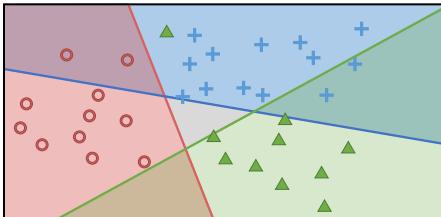
## Multiclass (more than 2) Classification

- **One-vs-rest** train separate binary classifiers for each class



## Multiclass (more than 2) Classification

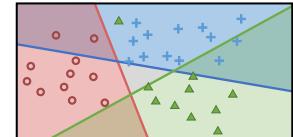
- **One-vs-rest** train separate binary classifiers for each class



## Multiclass (more than 2) Classification

- **One-vs-rest** train separate binary classifiers for each class

- Class with highest confidence wins
- Need to address class imbalance issue



- **Soft-Max** multiclass classification

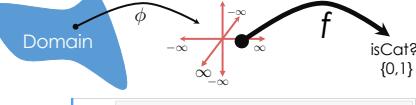
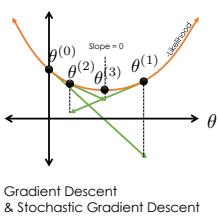
- **Soft-Max** multiclass classification

$$\mathbf{P}(Y = j | x) = \frac{\exp(x^T \theta^{(j)})}{\sum_{m=1}^k \exp(x^T \theta^{(m)})}$$

- Separate  $\theta^{(j)} \in \mathbb{R}^p$  for each class
- Trained using gradient descent methods
- Generates calibrated probabilities

Python Demo!

## Summary of Lecture



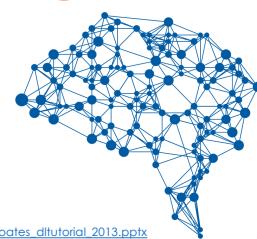
```
In [28]: def gradient_descent(X, Y, theta0, gradient_function, max_iter = 1000000, epsilon=0.0001):
    theta = theta0
    for t in range(1, max_iter):
        grad = gradient_function(theta, X, Y)
        theta = theta - gradient_function(theta, X, Y) * learning_rate
        # Track Convergence
        if np.linalg.norm(grad) < epsilon:
            return theta
    return theta
```

Logistic Regression

Multiclass Classification  
• One-vs-Rest  
• Softmax

## Deep Learning Overview

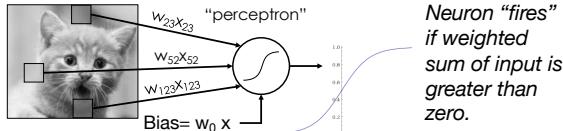
Bonus Material not on Exams!



Borrowed heavily from excellent talks by:  
• Adam Coates: [http://ai.stanford.edu/~acoates/coates\\_dltutorial\\_2013.pptx](http://ai.stanford.edu/~acoates/coates_dltutorial_2013.pptx)  
• Fei-Fei Li and Andrej Karpathy: <http://cs231n.stanford.edu/syllabus.html>

## Logistic Regression as a “Neuron”

➤ Consider the simple function family:  $\sigma(u) = \frac{1}{1 + \exp(-u)}$

$$f_w(x) = \sigma(w^T x) = \sigma\left(\sum_{j=1}^d w_j x_j\right) = P(y = 1 | x)$$


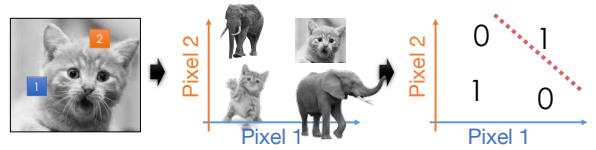
## Logistic Regression: Strengths and Limitations

➤ Widely used machine learning technique

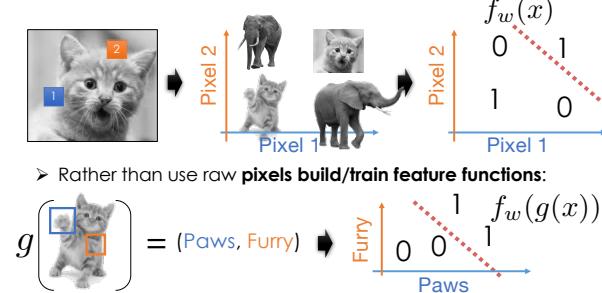
- convex → efficient to learn
- easy to interpret model weights
- works well given good features

➤ Limitations:

- Restricted to linear relationships → sensitive to choice of features



## Feature Engineering



## Composition Linear Models and Nonlinearities

$$\begin{matrix} d \\ k \end{matrix} \quad W^0 \quad \begin{matrix} d \\ k \\ 1 \end{matrix} \quad = \quad k \quad z \quad \begin{matrix} d \\ 1 \\ 1 \end{matrix} \quad \rightarrow \quad \sigma \begin{pmatrix} k \\ z \\ 1 \end{pmatrix} = \begin{matrix} h_1 \\ h_2 \\ \vdots \\ h_k \\ 1 \end{matrix}$$

## Composition Linear Models and Nonlinearities

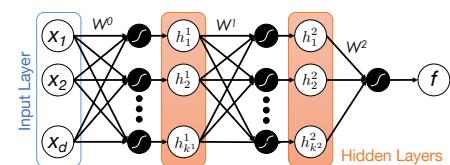
$$2 \quad W^1 \quad \begin{matrix} k \\ k \\ 1 \end{matrix} \quad = \quad 2 \quad z \quad \begin{matrix} 2 \\ 1 \end{matrix} \quad \rightarrow \quad \sigma \begin{pmatrix} 2 \\ z \\ 1 \end{pmatrix} = \begin{matrix} h_1 \\ h_2 \\ \vdots \\ h_k \\ 1 \end{matrix} \quad 2$$

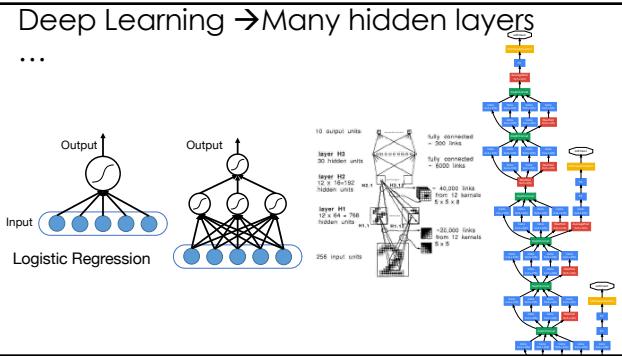
## Neural Networks

➤ Composing “perceptrons”

$$x \rightarrow \sigma(W^0 x) \rightarrow h^1 \rightarrow \sigma(W^1 h^1) \rightarrow h^2 \rightarrow \sigma(W^2 h^2) \rightarrow f$$

$$y = f_{W^0, W^1, W^2}(x) = \sigma(W^2 \sigma(W^1 \sigma(W^0 x)))$$





### Interested in Deep Learning?

- RISE Lab Deep Learning Overview:
  - [https://ucbrise.github.io/cs294-rise-fa16/deep\\_learning.html](https://ucbrise.github.io/cs294-rise-fa16/deep_learning.html)
- [TensorFlow Python Tutorial](#)
- Stanford CS231 Labs
  - <http://cs231n.github.io/linear-classify/>
  - <http://cs231n.github.io/optimization-1/>
  - <http://cs231n.github.io/optimization-2/>