

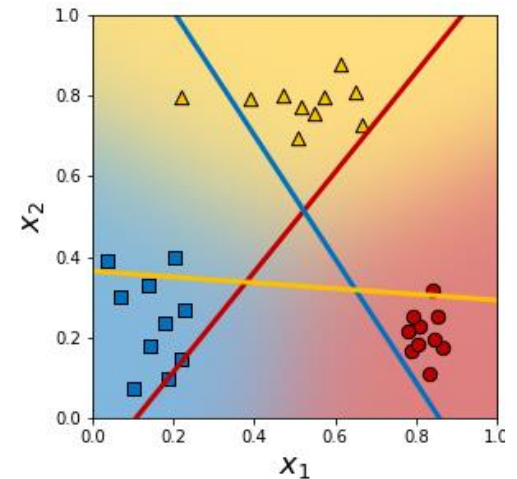
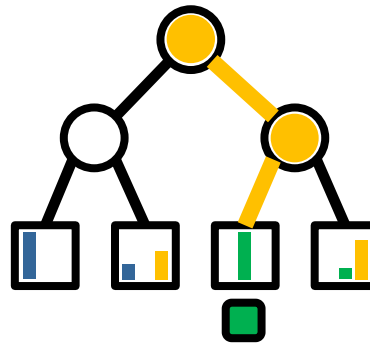
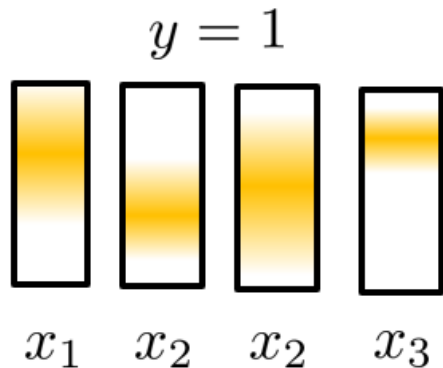
Photogrammetry & Robotics Lab

Machine Learning for Robotics and Computer Vision Tutorial

More on Classification

Jens Behley

Topics of this week lecture



- Classification models
 - Naïve Bayes (Generative Model)
 - Decision Tree (Discriminative Model)
 - Logistic/Softmax Regression (Discriminative Model)
- Optimization with Gradient Descent

Logistic Regression

- For **binary** classification $y = \{0, 1\}$, we want:

$$P(y = 0|\mathbf{x}) = 1 - P(y = 1|\mathbf{x})$$

- In **Logistic Regression** we define our model as:

$$\begin{aligned} P(y = 1|\mathbf{x}) &= \sigma(\theta^T \mathbf{x}) \\ &= \frac{1}{1 + \exp(-\theta^T \mathbf{x})} \end{aligned}$$

$$P(y = 0|\mathbf{x}) = 1 - P(y = 1|\mathbf{x})$$

- (We used again as in the Linear Regression: $\mathbf{x} := (1, \mathbf{x}^T)^T$)

Recap: Gradient of NLL

- For the gradient follows:

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \theta} &= \frac{1}{\partial \theta} \left(- \sum_{i=1}^N (\mathbf{1}\{y_i = 1\} - 1) \theta^T \mathbf{x} - \log(1 + \exp(-\theta^T \mathbf{x})) \right) \\&= - \sum_{i=1}^N (\mathbf{1}\{y_i = 1\} - 1) \mathbf{x} - \frac{1}{1 + \exp(-\theta^T \mathbf{x})} \exp(-\theta^T \mathbf{x}) (-\mathbf{x}) \\&= - \sum_{i=1}^N (\mathbf{1}\{y_i = 1\} - 1) \mathbf{x} + \underbrace{\frac{\exp(-\theta^T \mathbf{x})}{1 + \exp(-\theta^T \mathbf{x})}}_{1 - \sigma(\theta^T \mathbf{x})} \mathbf{x} \\&= \sum_{i=1}^N (\sigma(\theta^T \mathbf{x}) - \mathbf{1}\{y_i = 1\}) \mathbf{x} \\&= \sum_{i=1}^N (P(y_i = 1 | \mathbf{x}) - \mathbf{1}\{y_i = 1\}) \mathbf{x}\end{aligned}$$

- Problem:** Setting this to zero, no closed form solution!

How do you check the gradient?

Numerical Gradient

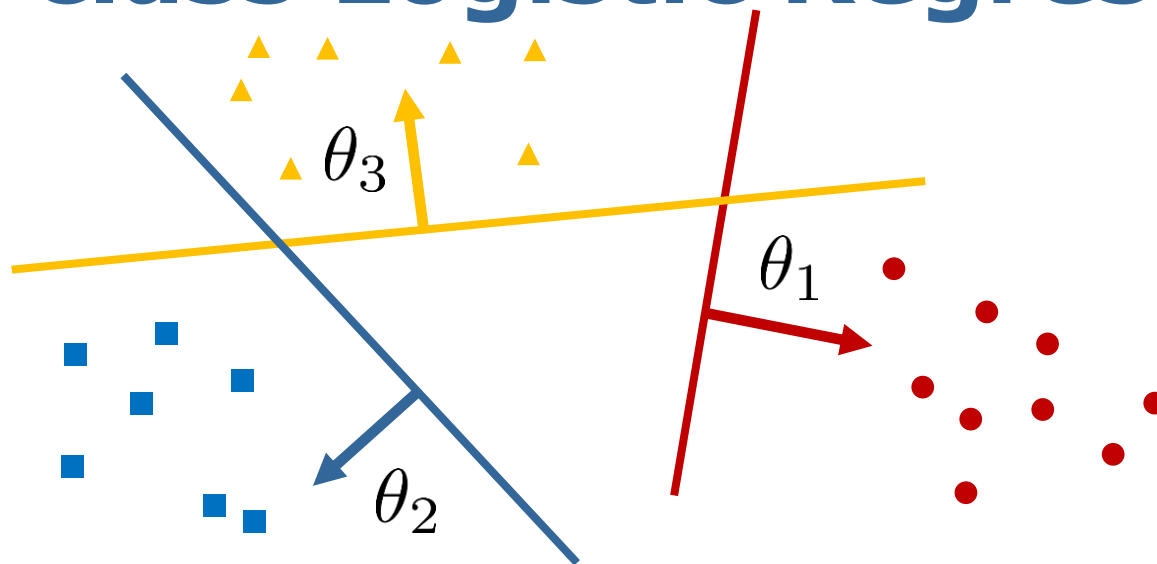
- Compute numerical gradient ($h = 1e-7$)

$$\frac{\hat{df}}{d\mathbf{x}_i} \approx \frac{f((x_0, \dots, x_i + h, \dots, x_D)^T) - f((x_0, \dots, x_i - h, \dots, x_D)^T)}{2 \cdot h}$$

- Relative error should be small (e.g., $1e-5$):

$$\text{rel_error} = \frac{\left| \frac{\hat{df}}{d\mathbf{x}_i} - \frac{df}{d\mathbf{x}_i} \right|}{\left| \frac{\hat{df}}{d\mathbf{x}_i} \right| + \left| \frac{df}{d\mathbf{x}_i} \right| + 1e-12}$$

Multi-class Logistic Regression



- **Idea:** class with largest $\mathbf{x}^T \theta_k$ should have highest confidence $P(y = k|\mathbf{x})$
- Want to find parameters such that distance is maximize for correct class
- How to turn “distances” into probability distribution?

Recap: Softmax

- Softmax function $\text{softmax} : \mathbb{R}^D \mapsto [0, 1]^D$ is given by

$$\text{softmax}(\mathbf{x}) = \mathbf{s}$$

with

$$s_i = \frac{\exp(x_i)}{\sum_{d=1}^D \exp(x_d)}$$

- Properties of softmax function:
 - $s_i \in [0, 1]$
 - $\sum_{i=1}^D s_i = 1$

Intuition of Softmax

- Some examples for intuition about the output of softmax function:
 - $\text{softmax}(10, 10, 10) = (1/3, 1/3, 1/3)$
 - $\text{softmax}(10, 11, 10) = (0.21, 0.58, 0.21)$
 - $\text{softmax}(10, 13, 10) = (0.045, 0.91, 0.045)$
 - $\text{softmax}(9, 11, 10) = (0.09, 0.67, 0.24)$

More on the softmax intuition

- $\text{softmax}(9, 11, 10) = (0.09, 0.67, 0.24)$
- $\text{softmax}(109, 111, 110) = ?$

More on the softmax intuition

- $\text{softmax}(9, 11, 10) = (0.09, 0.67, 0.24)$
- $\text{softmax}(109, 111, 110) = (0.09, 0.67, 0.24)$
- $\text{softmax}(1009, 1011, 1010) = ?$

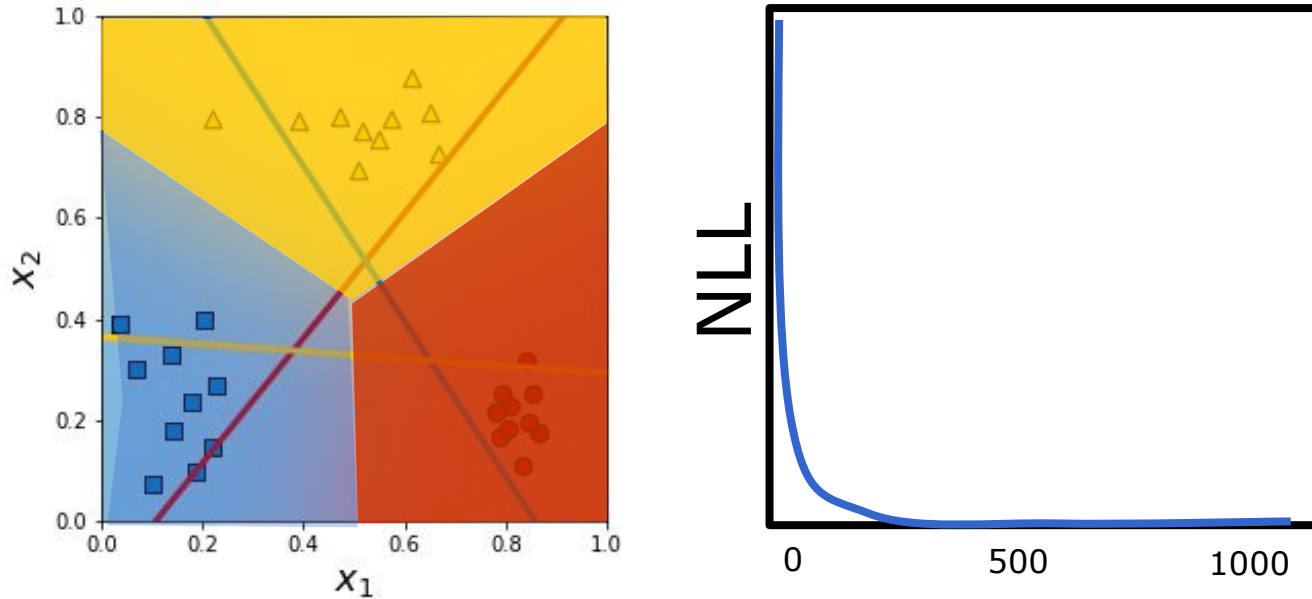
How to fix numerical overflow?

The softmax trick

$$\begin{aligned}\frac{\exp(a^{(i)})}{\exp(\sum_j a^{(j)})} &= \frac{\exp(a^{(i)}) \cdot \exp(z)}{\exp(\sum_j a^{(j)}) \cdot \exp(z)} \\ &= \frac{\exp(a^{(i)} + z)}{\exp(\sum_j a^{(j)} + z)}.\end{aligned}$$

- Thus, we can use $z = -\max_j a^{(j)}$ to get smaller arguments in the exponentiation.

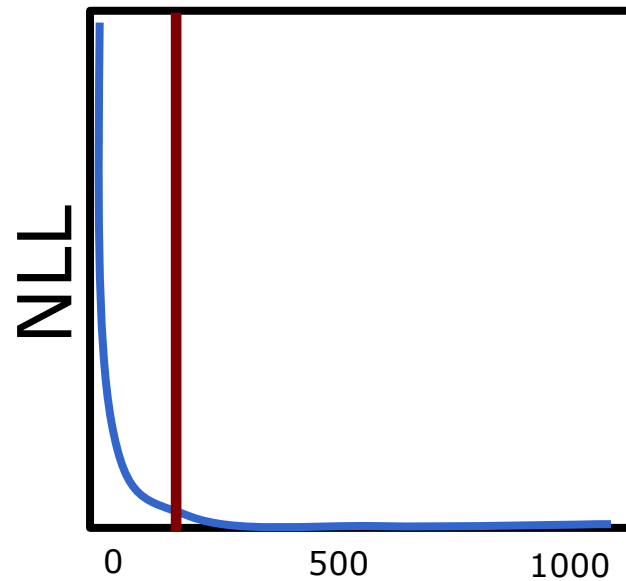
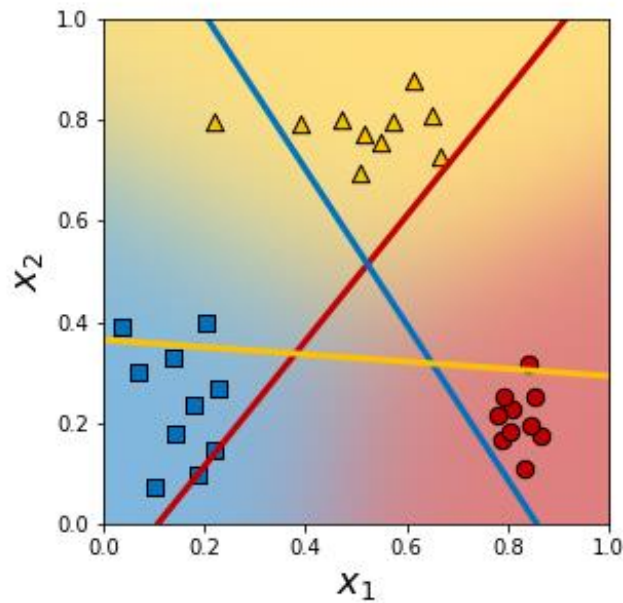
Training long enough...



- Train long enough and get “hard” boundaries

How to avoid overfitting?

Early Stopping



- Easy solution: stop early!
- Often used in NN training

See you next week!