

Deep Learning

A course about theory & practice



Regression vs. Classification

Marco Piastra

Function Approximation (a.k.a. Regression)

Regression

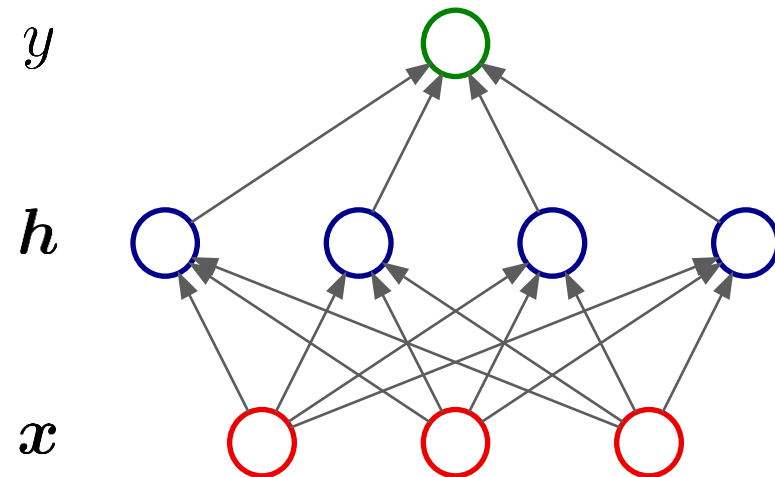
- **Function approximation**

$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d$$

Feed-forward neural network

$$\tilde{y} = \mathbf{w} \cdot g(\mathbf{W}\mathbf{x} + \mathbf{b}) + b$$

The Feed-forward neural network is used as an approximator of $f^*(\mathbf{x})$



Classification

Classification: Softmax

Classification

$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d, \quad y \in \{\text{class}_i\}_{i=1}^k$$

Feed-forward neural network with a *Softmax* layer

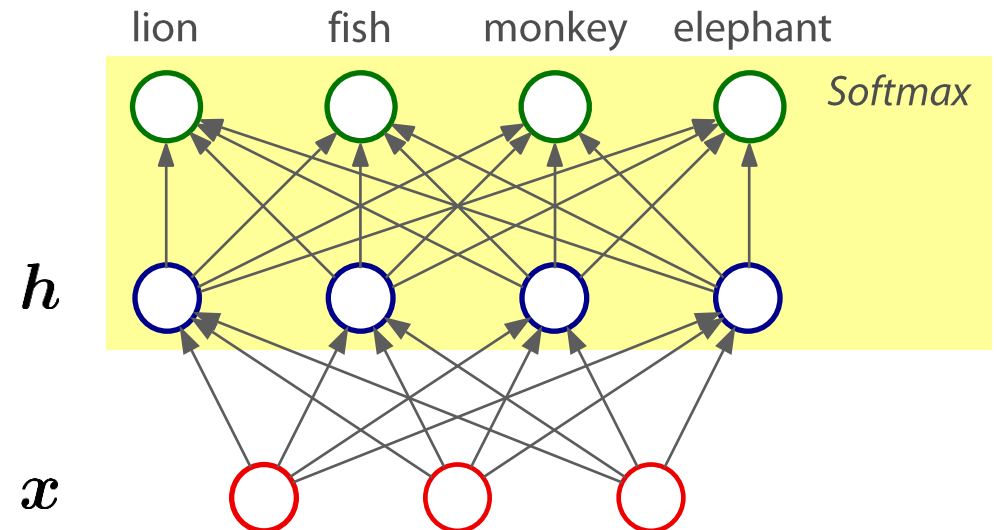
$$P(\tilde{y} = \text{class}_i | \mathbf{x}) := \frac{\exp(\mathbf{w}_i \cdot g(\mathbf{W}\mathbf{x} + \mathbf{b}) + b_i)}{\sum_{j=1}^k \exp(\mathbf{w}_j \cdot g(\mathbf{W}\mathbf{x} + \mathbf{b}) + b_j)}$$

From now on

$$P(\tilde{y} = \text{class}_i | \mathbf{x})$$

will be written as

$$P(\tilde{y} = i | \mathbf{x})$$



Classification: Softmax

■ Classification

$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d, \quad y \in \{\text{class}_i\}_{i=1}^k$$

The *Softmax* layer can be rewritten as:

$$P(\tilde{y} = \text{class}_i \mid \mathbf{h}) := \frac{\exp(\mathbf{w}_i \cdot \mathbf{h} + b_i)}{\sum_{j=1}^k \exp(\mathbf{w}_j \cdot \mathbf{h} + b_j)}$$

where, in this case: $\mathbf{h} := g(\mathbf{W}\mathbf{x} + \mathbf{b})$

(yet, more in general, \mathbf{h} can be anything)

Classification: Softmax

■ Softmax as a layer

The entire *Softmax* layer can be rewritten as:

$$P((\tilde{y} = i)_1^k | \mathbf{h}) := \frac{\exp(\mathbf{W}_S \mathbf{h} + \mathbf{b}_S)}{\sum \exp(\mathbf{W}_S \mathbf{h} + \mathbf{b}_S)}$$

Probability distribution
(a vector)

Sum of all components

where:

$$\mathbf{W}_S := \begin{bmatrix} - & \mathbf{w}_1 & - \\ & \vdots & \\ - & \mathbf{w}_k & - \end{bmatrix} \quad \mathbf{b}_S := \begin{bmatrix} b_1 \\ \vdots \\ b_k \end{bmatrix}$$

The vector $\mathbf{W}_S \mathbf{h} + \mathbf{b}_S$ is sometimes referred to as the **logit (vector)**

Classification: Softmax

■ Cross-entropy (in general)

P and Q are probability distributions on a discrete random variable $y \in \{1, \dots, k\}$

$$H(Q, P) := - \sum_{j=1}^k Q(y = j) \log P(\tilde{y} = j)$$

■ A loss function for Softmax

Q describes the 'true' class, i.e., the one in the dataset

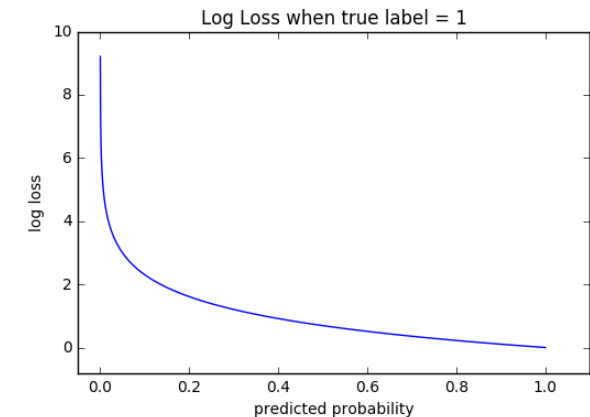
$$\hat{y}_j^{(i)} := \begin{cases} 1 & \text{if } y^{(i)} = j \\ 0 & \text{otherwise} \end{cases}$$

while P is the output of the Softmax layer

$$P(\tilde{y} = j | \mathbf{h})$$

Hence, the loss is:

$$\begin{aligned} L(\mathbf{h}^{(i)}, y^{(i)}) &:= - \sum_{j=1}^k \hat{y}_j^{(i)} \log P(\tilde{y} = j | \mathbf{h}^{(i)}) \\ &= - \log P(\tilde{y} = y^{(i)} | \mathbf{h}^{(i)}) \end{aligned}$$



Classification: Softmax

■ Cross-entropy for Softmax

$$L(\mathbf{h}^{(i)}, \mathbf{y}^{(i)}) := - \sum_{j=1}^k \delta(\mathbf{y}^{(i)} = j) \log P(\tilde{y} = j | \mathbf{h}^{(i)})$$

Expressing the loss function in vector form:

$$\mathbf{y} := \begin{bmatrix} y_1 \\ \vdots \\ y_k \end{bmatrix}, \quad y_j := \delta(\mathbf{y} = j) \quad \text{'one hot' representation} \quad \mathbf{p} := \begin{bmatrix} p_1 \\ \vdots \\ p_k \end{bmatrix}, \quad p_j := P(\tilde{y} = j | \mathbf{h})$$

$$L(\mathbf{h}^{(i)}, \mathbf{y}^{(i)}) = - \mathbf{y}^{(i)} \cdot \log(\mathbf{p}^{(i)})$$

which implies that also the dataset has to be transformed in the 'one hot' representation

$$D := \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N \quad \Longrightarrow \quad D := \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$$

Classification: Softmax

■ Gradient of Softmax (layerwise)

$$L(D) = \sum_{i=1}^N L(\mathbf{h}^{(i)}, \mathbf{y}^{(i)}) = - \sum_{i=1}^N \mathbf{y}^{(i)} \cdot \log(\mathbf{p}^{(i)})$$

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(D) &= \frac{\partial}{\partial \boldsymbol{\vartheta}} \sum_{i=1}^N L(\mathbf{h}^{(i)}, \mathbf{y}^{(i)}) = - \frac{\partial}{\partial \boldsymbol{\vartheta}} \sum_{i=1}^N \mathbf{y}^{(i)} \cdot \log(\mathbf{p}^{(i)}) \\ &= - \sum_{i=1}^N \mathbf{y}^{(i)} \cdot \frac{\partial}{\partial \boldsymbol{\vartheta}} \log(\mathbf{p}^{(i)}) \end{aligned}$$

This is a matrix

$$\frac{\partial}{\partial \boldsymbol{\vartheta}} \log(\mathbf{p}) = \begin{bmatrix} \frac{\partial}{\partial \vartheta_1} \log(p_1) & \dots & \frac{\partial}{\partial \vartheta_d} \log(p_1) \\ \vdots & \ddots & \vdots \\ \frac{\partial}{\partial \vartheta_1} \log(p_k) & \dots & \frac{\partial}{\partial \vartheta_d} \log(p_k) \end{bmatrix} = \begin{bmatrix} - & \frac{\partial}{\partial \boldsymbol{\vartheta}} \log(p_1) & - \\ & \vdots & \\ - & \frac{\partial}{\partial \boldsymbol{\vartheta}} \log(p_k) & - \end{bmatrix}$$

Classification: Softmax

- **Gradient of Softmax** (layerwise)

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\vartheta}} \log(p_j) &= \frac{\partial}{\partial \boldsymbol{\vartheta}} \log P(\tilde{y} = j \mid \mathbf{h}) \\ &= \frac{\partial}{\partial \boldsymbol{\vartheta}} \log \frac{\exp(\mathbf{w}_j \cdot \mathbf{h} + b_j)}{\sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l)} \\ &= \frac{\partial}{\partial \boldsymbol{\vartheta}} \left(\log \exp(\mathbf{w}_j \cdot \mathbf{h} + b_j) - \log \sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l) \right) \\ &= \frac{\partial}{\partial \boldsymbol{\vartheta}} (\mathbf{w}_j \cdot \mathbf{h} + b_j) - \frac{\partial}{\partial \boldsymbol{\vartheta}} \log \sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l)\end{aligned}$$

Classification: Softmax

■ Gradient of Softmax (layerwise)

$$\frac{\partial}{\partial \vartheta} \log(p_j) = \frac{\partial}{\partial \vartheta} (\mathbf{w}_j \cdot \mathbf{h} + b_j) - \frac{\partial}{\partial \vartheta} \log \sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l)$$

Case 1: $\vartheta = w_r$ or $\vartheta = b_r$

Case 2: $\mathbf{h}(\vartheta)$ i.e. ϑ is a generic parameter on which \mathbf{h} depends

$$\frac{\partial \mathbf{h}^{[i]}}{\partial \vartheta^{[i]}}$$
$$\frac{\partial \mathbf{h}^{[i]}}{\partial \vartheta^{[j]}}, \quad j < i$$

Let's compute the two contributions separately

$$\frac{\partial}{\partial \vartheta} (\mathbf{w}_j \cdot \mathbf{h} + b_j)$$
$$\frac{\partial}{\partial \vartheta} \log \sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l)$$

Classification: Softmax

■ Gradient of Softmax (layerwise)

$$\frac{\partial}{\partial \vartheta} (\mathbf{w}_j \cdot \mathbf{h} + b_j)$$

Case 1: $\vartheta = w_r$ or $\vartheta = b_r$

$$\frac{\partial}{\partial w_r} (\mathbf{w}_j \cdot \mathbf{h} + b_j) = \begin{cases} \mathbf{0} & \text{if } r \neq j \\ \mathbf{h} & \text{otherwise} \end{cases}$$

$$\frac{\partial}{\partial b_r} (\mathbf{w}_j \cdot \mathbf{h} + b_j) = \begin{cases} 0 & \text{if } r \neq j \\ 1 & \text{otherwise} \end{cases}$$

Case 2: $\mathbf{h}(\vartheta)$ i.e. ϑ is a generic parameter on which \mathbf{h} depends

$$\frac{\partial}{\partial \vartheta} (\mathbf{w}_j \cdot \mathbf{h} + b_j) = \mathbf{w}_j \cdot \frac{\partial}{\partial \vartheta} \mathbf{h}$$

Classification: Softmax

■ Gradient of Softmax (layerwise)

$$\frac{\partial}{\partial \vartheta} \log \sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l)$$

Case 1: $\vartheta = w_r$ or $\vartheta = b_r$

$$\begin{aligned} \frac{\partial}{\partial w_r} \log \sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l) &= \\ &= \frac{1}{\sum_{m=1}^k \exp(\mathbf{w}_m \cdot \mathbf{h} + b_m)} \frac{\partial}{\partial w_r} \sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l) \\ &= \frac{1}{\sum_{m=1}^k \exp(\mathbf{w}_m \cdot \mathbf{h} + b_m)} \sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l) \frac{\partial}{\partial w_r} (\mathbf{w}_l \cdot \mathbf{h} + b_l) \\ &= \frac{\exp(\mathbf{w}_r \cdot \mathbf{h} + b_r)}{\sum_{m=1}^k \exp(\mathbf{w}_m \cdot \mathbf{h} + b_m)} \mathbf{h} = p_r \mathbf{h} \end{aligned}$$

Classification: Softmax

■ Gradient of Softmax (layerwise)

$$\frac{\partial}{\partial \vartheta} \log \sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l)$$

Case 1: $\vartheta = w_r$ or $\vartheta = b_r$

$$\begin{aligned} \frac{\partial}{\partial b_r} \log \sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l) &= \\ &= \frac{1}{\sum_{m=1}^k \exp(\mathbf{w}_m \cdot \mathbf{h} + b_m)} \frac{\partial}{\partial b_r} \sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l) \\ &= \frac{1}{\sum_{m=1}^k \exp(\mathbf{w}_m \cdot \mathbf{h} + b_m)} \sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l) \frac{\partial}{\partial b_r} (\mathbf{w}_l \cdot \mathbf{h} + b_l) \\ &= \frac{\exp(\mathbf{w}_r \cdot \mathbf{h} + b_r)}{\sum_{m=1}^k \exp(\mathbf{w}_m \cdot \mathbf{h} + b_m)} = p_r \end{aligned}$$

Classification: Softmax

■ Gradient of Softmax (layerwise)

$$\frac{\partial}{\partial \vartheta} \log \sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l)$$

Case 2: $\mathbf{h}(\vartheta)$ i.e. ϑ is a generic parameter on which \mathbf{h} depends

$$\begin{aligned} \frac{\partial}{\partial \vartheta} \log \sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l) &= \\ &= \frac{1}{\sum_{m=1}^k \exp(\mathbf{w}_m \cdot \mathbf{h} + b_m)} \frac{\partial}{\partial \vartheta} \sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l) \\ &= \frac{1}{\sum_{m=1}^k \exp(\mathbf{w}_m \cdot \mathbf{h} + b_m)} \sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l) \frac{\partial}{\partial \vartheta} (\mathbf{w}_l \cdot \mathbf{h} + b_l) \\ &= \sum_{l=1}^k \frac{\exp(\mathbf{w}_l \cdot \mathbf{h} + b_l)}{\sum_{m=1}^k \exp(\mathbf{w}_m \cdot \mathbf{h} + b_m)} \mathbf{w}_l^T \frac{\partial}{\partial \vartheta} \mathbf{h} = \left(\sum_{l=1}^k p_l \mathbf{w}_l^T \right) \frac{\partial}{\partial \vartheta} \mathbf{h} \end{aligned}$$

Likelihood
(or the Mother of All Loss Functions)

MSE: Choosing the Noise Model

- **Function approximation**

$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d$$

Feed-forward neural network

$$\tilde{y} = \tilde{f}(\mathbf{x}; \boldsymbol{\vartheta}) = \mathbf{w} \cdot g(\mathbf{W}\mathbf{x} + \mathbf{b}) + b$$

Assume that the dataset is *noisy*

$$D := \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$$

Gaussian noise, constant variance

$$\text{where } y^{(i)} \sim \mathcal{N}(\mu^{(i)}, \sigma^2), \quad \mu^{(i)} = f^*(\mathbf{x}^{(i)}), \quad \forall i$$

Then, the approximator function can be interpreted as:

$$\mu = \tilde{f}(\mathbf{x}; \boldsymbol{\vartheta})$$

MSE: Choosing the Noise Model

■ Likelihood of the dataset

Assuming all data items $(\mathbf{x}^{(i)}, y^{(i)})$ are **Independent and Identically Distributed (I.I.D.)**

Joint Probability
of the dataset,
given $\boldsymbol{\vartheta}$

$$\begin{aligned}\mathcal{L}(D; \boldsymbol{\vartheta}) &= \prod_{i=1}^N P(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\vartheta}) \\ &= \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y^{(i)} - \mu^{(i)})^2}{2\sigma^2}\right)\end{aligned}$$

Gaussian noise, constant variance

$$\begin{aligned}\log \mathcal{L}(D; \boldsymbol{\vartheta}) &= \sum_{i=1}^N \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y^{(i)} - \mu^{(i)})^2}{2\sigma^2}\right) \right] \\ &= \sum_{i=1}^N \left(\log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \right] + \log \left[\exp\left(-\frac{(y^{(i)} - \mu^{(i)})^2}{2\sigma^2}\right) \right] \right) \\ &= \sum_{i=1}^N \left(-\frac{1}{2} \log(2\pi\sigma^2) - \frac{(y^{(i)} - \mu^{(i)})^2}{2\sigma^2} \right)\end{aligned}$$

MSE: Choosing the Noise Model

▪ Likelihood of the dataset

Assuming all data items $(\mathbf{x}^{(i)}, y^{(i)})$ are **Independent and Identically Distributed (I.I.D.)**

$$\log \mathcal{L}(D; \boldsymbol{\vartheta}) = \sum_{i=1}^N \left(-\frac{1}{2} \log(2\pi\sigma^2) - \frac{(y^{(i)} - \mu^{(i)})^2}{2\sigma^2} \right)$$

Actually, σ^2 is constant and it is not involved in the maximization:

$$\operatorname{argmax}_{\boldsymbol{\vartheta}} \log \mathcal{L}(D; \boldsymbol{\vartheta}) = - \sum_{i=1}^N (y^{(i)} - \mu^{(i)})^2$$

By convention, *minimization* (i.e., gradient descent) is preferred:

$$\operatorname{argmin}_{\boldsymbol{\vartheta}} \text{MSE}(D; \boldsymbol{\vartheta}) = \sum_{i=1}^N (y^{(i)} - \mu^{(i)})^2$$

Moral: *Mean Squared Error* (MSE) is the right loss function when the **noise model** is Gaussian with constant σ^2

Softmax: End-to-End Likelihood

■ Classification

$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d, \quad y \in \{\text{class}_i\}_{i=1}^k$$

The Softmax layer (as a function of the *logit vector*):

$$P(\tilde{y} = y_c \mid \mathbf{z}; \boldsymbol{\vartheta}) := \frac{\text{exp}(z_c)}{\sum_{j=1}^k \text{exp}(z_j)}$$

class c \ / *logit vector*

Softmax: End-to-End Likelihood

▪ Likelihood of the dataset

Assuming all data items $(\mathbf{x}^{(i)}, y^{(i)})$ are **Independent and Identically Distributed (I.I.D.)**

Joint Probability
of the dataset,
given $\boldsymbol{\vartheta}$

$$\begin{aligned}\mathcal{L}(D; \boldsymbol{\vartheta}) &= \prod_{i=1}^N P(\tilde{y}^{(i)} = y^{(i)} \mid \mathbf{z}^{(i)}; \boldsymbol{\vartheta}) \\ &= \prod_{i=1}^N \frac{\exp(z_c^{(i)})}{\sum_{j=1}^K \exp(z_j^{(i)})}\end{aligned}$$

Define an *indicator* (i.e., one-hot encoding)

$$\hat{y}_c^{(i)} := \mathbb{I}(y^{(i)} = c) = \begin{cases} 1 & \text{if } y^{(i)} = c \\ 0 & \text{otherwise} \end{cases}$$

The likelihood becomes:

$$\mathcal{L}(D; \boldsymbol{\vartheta}) = \prod_{i=1}^N \prod_{c=1}^K \left[\frac{\exp(z_c^{(i)})}{\sum_{j=1}^K \exp(z_j^{(i)})} \right]^{\hat{y}_c^{(i)}}$$

Softmax: End-to-End Likelihood

▪ Likelihood of the dataset

Assuming all data items $(\mathbf{x}^{(i)}, y^{(i)})$ are **Independent and Identically Distributed (I.I.D.)**

$$\begin{aligned}\log \mathcal{L}(D; \boldsymbol{\vartheta}) &= \sum_{i=1}^N \sum_{c=1}^K \hat{y}_c^{(i)} \log \left[\frac{\exp(z_c^{(i)})}{\sum_{j=1}^K \exp(z_j^{(i)})} \right] \\ &= \sum_{i=1}^N \sum_{c=1}^K \hat{y}_c^{(i)} \log P(\tilde{y}^{(i)} = y_c \mid \mathbf{z}^{(i)}; \boldsymbol{\vartheta}) \\ &= \sum_{i=1}^N \log P(\tilde{y}^{(i)} = y^{(i)} \mid \mathbf{z}^{(i)}; \boldsymbol{\vartheta})\end{aligned}$$

the log of a probability is always negative or zero

By convention, *minimization* (i.e., gradient descent) is preferred:

$$\operatorname{argmin}_{\boldsymbol{\vartheta}} \text{CE}(D; \boldsymbol{\vartheta}) = - \sum_{i=1}^N \log P(\tilde{y}^{(i)} = y^{(i)} \mid \mathbf{z}^{(i)}; \boldsymbol{\vartheta})$$

Moral: *Cross Entropy* (CE) is the right loss function when the class attribution is **uncertain (noisy)**

Regularization

Beyond Likelihood: Regularization

- A penalty on complexity

$$L_{total} = L_{data} + \lambda R(\boldsymbol{\vartheta})$$

coefficient λ
Penalty on complexity $R(\boldsymbol{\vartheta})$

- Norm-based regularization

L_2 Regularization (Weight Decay)

$$R(\boldsymbol{\vartheta}) = \frac{1}{2} \|\boldsymbol{\vartheta}\|_2^2 = \frac{1}{2} \boldsymbol{\vartheta}^T \boldsymbol{\vartheta} = \frac{1}{2} (\vartheta_1^2 + \vartheta_2^2 + \dots + \vartheta_n^2)$$

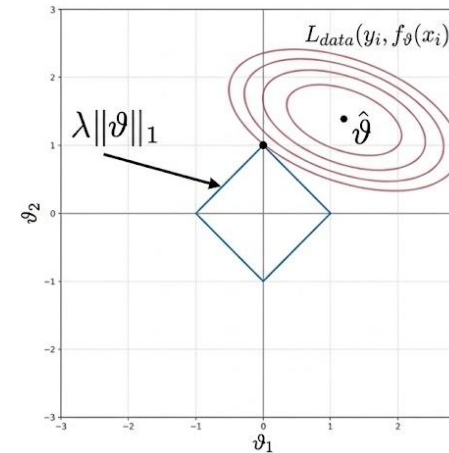
Geometrically, it pushes the solution toward a sphere centered at the origin
Smooths the function by preventing any single feature from dominating

L_1 Regularization (Lasso)

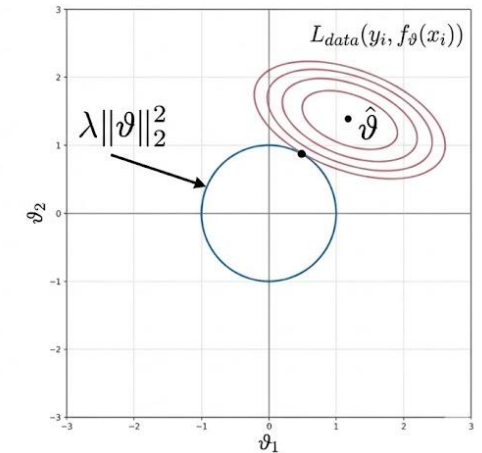
$$R(\boldsymbol{\vartheta}) = \|\boldsymbol{\vartheta}\|_1 = \sum_{i=1}^n |\vartheta_i| = |\vartheta_1| + |\vartheta_2| + \dots + |\vartheta_n|$$

Promotes sparsity (forces many weights to exactly zero). Useful for feature selection.

L1 Regularization



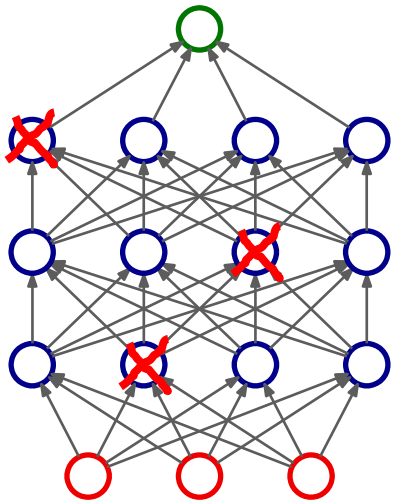
L2 Regularization



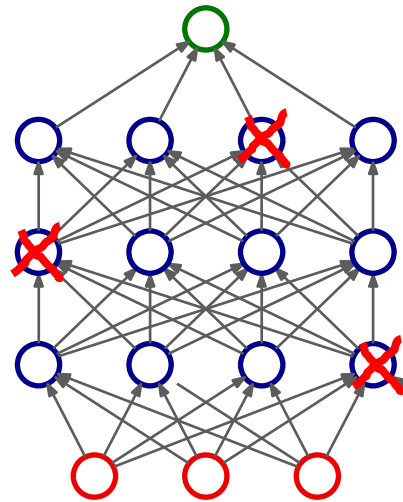
Stochastic Regularization: Dropout

■ Knocking-out at random

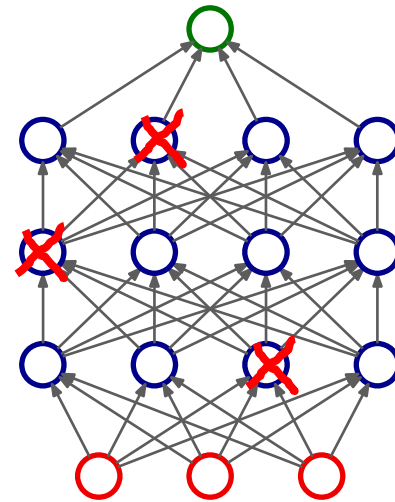
For each mini-batch, a small percentage of 'units' is de-activated



Training: mini-batch 1

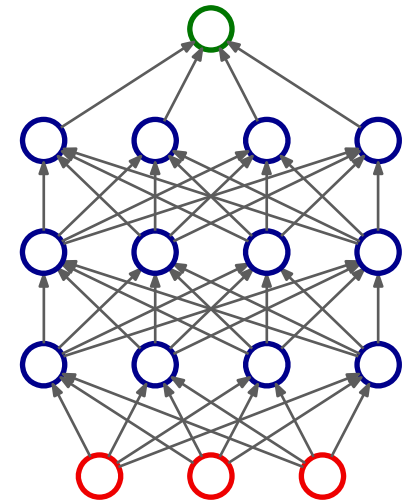


Training: mini-batch 2



Training: mini-batch 3

At inference time,
dropout is not active



Inference Time

Instead of penalizing the values of weights, it promotes the reliability of the network
It prevents **co-adaptation**: each unit must learn robust, independent features