**6.1 Training Neural Networks**

*Machine Learning 1: Foundations*

Marius Kloft  *(TUK)*

# Recap

Artificial neural networks (ANN)

- ▶ Key advantage over SVM, logistic regression, and friends: can **learn a good representation** of the data,

$$\min_{b\in\mathbb{R},\mathbf{w}\in\mathbb{R}^d} \quad \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_{i=1}^{n}\log\left(0, 1+\exp(-y_i(\mathbf{w}^\top \mathbf{x}_i + b))\right).$$

# Recap

Artificial neural networks (ANN)

▶ Key advantage over SVM, logistic regression, and friends: can **learn a good representation** of the data,

$$\min_{b\in\mathbb{R},\mathbf{w}\in\mathbb{R}^d} \quad \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_{i=1}^{n}\log\left(0, 1+\exp(-y_i(\mathbf{w}^\top\phi(\mathbf{x}_i)+b))\right).$$

# Recap

Artificial neural networks (ANN)

▶ Key advantage over SVM, logistic regression, and friends: can **learn a good representation** of the data,

$$\min_{b\in\mathbb{R},\mathbf{w}\in\mathbb{R}^d,\phi} \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_{i=1}^{n}\log\left(0, 1+\exp(-y_i(\mathbf{w}^\top\phi(\mathbf{x}_i)+b))\right).$$

# Recap

Artificial neural networks (ANN)

- ▶ Key advantage over SVM, logistic regression, and friends: can **learn a good representation** of the data,

$$\min_{b\in\mathbb{R},\mathbf{w}\in\mathbb{R}^d,\phi} \quad \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_{i=1}^{n}\log\left(0, 1+\exp(-y_i(\mathbf{w}^\top\phi(\mathbf{x}_i)+b))\right).$$

Need to restrict search space of $\phi$!
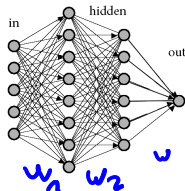
# Recap

Artificial neural networks (ANN)

▶ Key advantage over SVM, logistic regression, and friends: can **learn a good representation** of the data,

$$\min_{b\in\mathbb{R},\mathbf{w}\in\mathbb{R}^d,\phi} \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_{i=1}^{n}\log\Big(0, 1+\exp(-y_i(\mathbf{w}^\top\phi(\mathbf{x}_i)+b))\Big).$$

Need to restrict search space of $\phi$!

Idea: design $\phi$ similar to our brain

▶ multiple neurons in multiple layers with feed-forward connections

▶ $\phi_W(\mathbf{x}_i) := \sigma\Big(W_{L-1}^\top\ldots\sigma(W_1^\top\cdot\mathbf{x}_i)\ldots\Big)$

▶ optimize over $W = \big(W_1,\ldots,W_L\big)$!
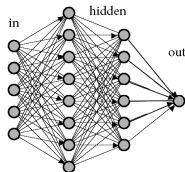
# Recap

Artificial neural networks (ANN)

▶ Key advantage over SVM, logistic regression, and friends: can **learn a good representation** of the data,

$$\min_{b\in\mathbb{R},\mathbf{w}\in\mathbb{R}^d,\phi} \quad \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_{i=1}^{n}\log\left(0, 1+\exp(-y_i(\mathbf{w}^\top\phi(\mathbf{x}_i)+b))\right).$$

> Need to restrict search space of $\phi$!

Idea: design $\phi$ similar to our brain

▶ multiple neurons in multiple layers with feed-forward connections

▶ $\phi_W(\mathbf{x}_i) := \sigma\left(W_{L-1}^\top \ldots \sigma(W_1^\top \cdot \mathbf{x}_i)\ldots\right)$

▶ optimize over $W = (W_1, \ldots, W_L)$!



> How to train ANNs?

# Contents of this Class

# How to Train (Deep) ANNs?

---

For the sake of simplicity, we focus on discussing how to train *fully connected ANNs* (not CNNs).

# How to Train (Deep) ANNs?

In the same way as we trained the SVM:
**(stochastic) gradient descent**!

For the sake of simplicity, we focus on discussing how to train *fully connected ANNs* (not CNNs).

# How to Train (Deep) ANNs?

> In the same way as we trained the SVM:
> **(stochastic) gradient descent**!

Recall the ANN optimization problem:

$$\min_{\mathbf{w}, W} \underbrace{\frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{2} \sum_{l=1}^{L} \|W_l\|_{\text{Fro}}^2 + C \sum_{i=1}^{n} \log \left( 1 + \exp \left( - y_i \mathbf{w}^\top \phi_W(\mathbf{x}_i) \right) \right)}_{=:F(\mathbf{w}, W)}$$

---

For the sake of simplicity, we focus on discussing how to train *fully connected ANNs* (not CNNs).

# How to Train (Deep) ANNs?

> In the same way as we trained the SVM:
> **(stochastic) gradient descent**!

Recall the ANN optimization problem:

$$\min_{\mathbf{w}, W} \underbrace{\frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{2} \sum_{l=1}^{L} \|W_l\|_{\text{Fro}}^2 + C \sum_{i=1}^{n} \log\left(1 + \exp\left(-y_i \mathbf{w}^\top \phi_W(\mathbf{x}_i)\right)\right)}_{=:F(\mathbf{w}, W)}$$

> How to compute the gradient of $F$?

---

For the sake of simplicity, we focus on discussing how to train *fully connected ANNs* (not CNNs).

# The Gradient of $F$ With Respect to $\boxed{\mathbf{w}}$ is Simple:

$$\nabla_{\mathbf{w}} F(\mathbf{w}, W) = \mathbf{w} + C \sum_{i=1}^{n} \nabla_{\mathbf{w}} \log \left( 1 + \exp \left( -y_i \mathbf{w}^{\top} \phi_W(\mathbf{x}_i) \right) \right)$$

$$\overset{(*)}{=} \mathbb{1} + c \sum_{i=1}^{n} \frac{-y_i \, \phi_\omega(x_i)}{\mathbb{1} + \exp \left( y_i \, w^{\top} \phi_\omega(x_i) \right)}$$

$$= \mathbb{1} - c \sum_{i=1}^{n} \frac{y_i \, \phi_\omega(x_i)}{\mathbb{1} + \exp \left( y_i \, w^{\top} \phi_\omega(x_i) \right)}$$

$$\overset{(*)}{\Big|} \, \log \left( \mathbb{1} + \exp(t) \right)$$

$$= \frac{\exp(t)}{\mathbb{1} + \exp(t)} = \frac{\mathbb{1}}{\mathbb{1} + \exp(-t)}$$

# The Gradient of $F$ With Respect to $\boxed{\mathbf{w}}$ is Simple:

$$\nabla_{\mathbf{w}} F(\mathbf{w}, W) = \mathbf{w} + C \sum_{i=1}^{n} \nabla_{\mathbf{w}} \log \left( 1 + \exp \left( -y_i \mathbf{w}^{\top} \phi_W(\mathbf{x}_i) \right) \right)$$

But how to compute the gradient of $F$ with respect to $\boxed{W}$?

# Gradient of $F$ With Respect to $\boxed{W = (W_1, \ldots, W_L)}$

Analogously, we have, for all $l = 1, \ldots, L$:

$$\nabla_{W_l} F(\mathbf{w}, W) = W_l + C \sum_{i=1}^{n} \nabla_{W_l} \log \Big( 1 + \exp \big( - y_i \mathbf{w}^\top \phi_W(\mathbf{x}_i) \big) \Big)$$

$$= W_l + C \sum_{i=1}^{n} \frac{- y_i \, \omega^\top \nabla_{\omega} \phi_{\omega}(x_i)}{1 + \exp \big( y_i \, \omega^\top \phi_{\omega}(x_i) \big)}$$

From now on, denote the $ij$th entry of $W_l$ by $w_{ijl}$.

# Gradient of $F$ With Respect to $\boxed{W = (W_1, \ldots, W_L)}$

Analogously, we have, for all $l = 1, \ldots, L$:

$$\nabla_{W_l} F(\mathbf{w}, W) = W_l + C \sum_{i=1}^n \nabla_{W_l} \log \left( 1 + \exp\left( - y_i \mathbf{w}^\top \phi_W(\mathbf{x}_i) \right) \right)$$

From now on, denote the $ij$th entry of $W_l$ by $w_{ijl}$.

Given a data point $\mathbf{x}$, how to compute $\boxed{\nabla_{w_{ijl}} \phi_W(\mathbf{x})}$ ?

# Computing $\nabla_{w_{ijl}}\phi_W(\mathbf{x})$

We have:

$$\nabla_{w_{ijl}}\phi_W(\mathbf{x}) = \nabla_{w_{ijl}}\,\sigma\Big(\,W_L^\top\sigma\big(\ldots\sigma(\underbrace{\underbrace{\underbrace{\underbrace{W_1^\top\mathbf{v}_0}_{=\mathbf{u}_1}}_{=\mathbf{v}_1}}_{\mathbf{u}_L}}_{\mathbf{v}_L})\ldots\big)\Big).$$

# Computing $\nabla_{w_{ijl}} \phi_W(\mathbf{x})$

We have:

$$\nabla_{w_{ijl}} \phi_W(\mathbf{x}) = \nabla_{w_{ijl}} \, \sigma \Big( \, W_L^\top \sigma \big( \ldots \sigma (\underbrace{W_1^\top \mathbf{v}_0}_{=\mathbf{u}_1}) \ldots \big) \Big).$$

$$\underbrace{\phantom{W_1^\top \mathbf{v}_0}}_{=\mathbf{v}_1}$$

$$\underbrace{\phantom{W_L^\top \sigma \big( \ldots \sigma (W_1^\top \mathbf{v}_0)}}_{\mathbf{u}_L}$$

$$\underbrace{\phantom{\sigma \big( W_L^\top \sigma \big( \ldots \sigma (W_1^\top \mathbf{v}_0) \big)}}_{\mathbf{v}_L}$$

Need to compute a gradient of a **nested** function!

# Computing $\nabla_{w_{ijl}} \phi_W(\mathbf{x})$

We have:

$$\nabla_{w_{ijl}} \phi_W(\mathbf{x}) = \nabla_{w_{ijl}} \sigma\Big( W_L^\top \sigma\big(\ldots \sigma(\underbrace{W_1^\top \mathbf{v}_0}_{=\mathbf{u}_1})\ldots\big)\Big).$$

$$\underbrace{\phantom{W_1^\top \mathbf{v}_0}}_{=\mathbf{v}_1}$$

$$\underbrace{\phantom{W_L^\top \sigma\big(\ldots \sigma(W_1^\top \mathbf{v}_0)\ldots\big)}}_{\mathbf{u}_L}$$

$$\underbrace{\phantom{\sigma\Big( W_L^\top \sigma\big(\ldots \sigma(W_1^\top \mathbf{v}_0)\ldots\big)\Big)}}_{\mathbf{v}_L}$$

> Need to compute a gradient of a **nested** function!

## Idea: Chain rule

$$\nabla_{w_{ijl}} \phi_W(\mathbf{x}) = \frac{\partial \mathbf{v}_L}{\partial w_{ijl}} = \frac{\partial \mathbf{v}_L}{\partial \mathbf{u}_L} \cdot \frac{\partial \mathbf{u}_L}{\partial \mathbf{v}_{L-1}} \cdot \frac{\partial \mathbf{v}_{L-1}}{\partial \mathbf{u}_{L-1}} \cdots \frac{\partial \mathbf{u}_{l+1}}{\partial \mathbf{v}_l} \cdot \frac{\partial \mathbf{v}_l}{\partial \mathbf{u}_l} \cdot \frac{\partial \mathbf{u}_l}{\partial w_{ijl}}$$

① ② ① ② ① ③

# Three Terms Occur by the Chain Rule:

For all $l = 1, \ldots, L$:

1. $\frac{\partial \mathbf{v}_l}{\partial \mathbf{u}_l}$

2. $\frac{\partial \mathbf{u}_l}{\partial \mathbf{v}_{l-1}}$

3. $\frac{\partial \mathbf{u}_l}{\partial w_{ijl}}$

We need to compute all of them!

# First Term

We compute the first term as:

① $\quad \dfrac{\partial \mathbf{v}_l}{\partial \mathbf{u}_l} = \dfrac{\delta \max(0, u_\ell)}{\delta u_\ell}$

$$= \begin{cases} 1 & \text{if } u_\ell \geq 0 \\ 0 & \text{else wise} \end{cases}$$

$$= \Theta(u_\ell)$$

# Second Term

We compute the second term as:

② $\quad \dfrac{\partial \mathbf{u}_l}{\partial \mathbf{v}_{l-1}} = \dfrac{\delta \, w_\ell^\top v_{\ell-1}}{\delta v_{\ell-1}} = w_\ell^\top$

# Third Term

We compute the third term as:

③ $\quad \dfrac{\partial \mathbf{u}_l}{\partial w_{ijl}} = \dfrac{\delta\left( w_\ell^T v_{\ell-1} \right)}{\delta w_{ij\ell}} = \dfrac{\delta\left( \sum_k w_{k k' \ell} v_{k, \ell-1} \right)_{k'}}{\delta w_{ij\ell}}$

$\quad\quad = v_{i, \ell-1} \; e_j$

Notation:

- $v_{k, \ell-1}$ .... kth entry of $v_{\ell-1}$
- $e_j$ .... unit vector with
  $\quad\quad$ ↗ in jth component

# Putting Things Together

Our chain rule formula from Slide 9 thus translates into:

$$\nabla_{w_{ijl}} \phi_W(\mathbf{x}) = \frac{\partial \mathbf{v}_L}{\partial \mathbf{u}_L} \cdot \frac{\partial \mathbf{u}_L}{\partial \mathbf{v}_{L-1}} \cdot \frac{\partial \mathbf{v}_{L-1}}{\partial \mathbf{u}_{L-1}} \cdots \frac{\partial \mathbf{u}_{l+1}}{\partial \mathbf{v}_l} \cdot \frac{\partial \mathbf{v}_l}{\partial \mathbf{u}_l} \cdot \frac{\partial \mathbf{u}_l}{\partial w_{ijl}}$$

$$= \Theta(\mathbf{u}_L) W_L^\top \Theta(\mathbf{u}_{L-1}) \cdots W_{l+1}^\top \Theta(\mathbf{u}_l) v_{i,l-1} \mathbf{e}_j$$

How to code up the computation of

$$\nabla_{w_{ijl}} \phi_W(\mathbf{x}) \qquad \forall i, j, l$$

in an efficient algorithm?

# Backpropagation Algorithm

Given an input **x**, we first compute all variables $\mathbf{u}_l$ and $\mathbf{v}_l$:

## Forward propagation

1: initialize $\mathbf{v}_0 := \mathbf{x}$
2: **for** $l = 1 : (L-1)$ **do**
3: $\quad \mathbf{u}_l := W_l^\top \mathbf{v}_{l-1}$
4: $\quad \mathbf{v}_l := \Theta(\mathbf{u}_l)$
5: **end for**

$x \rightarrow u_1 \rightarrow v_1 \rightarrow u_2 \rightarrow v_2 \rightarrow \ldots$
$\rightarrow v_L$

Then, we compute the gradient via the chain rule:

## Backward propagation

1: initialize $\delta_L := \Theta(\mathbf{u}_L)$
2: $\nabla_{w_{ijL}} \phi_W(\mathbf{x}) := \delta_L v_{i,L-1} \mathbf{e}_j \qquad \forall i, j$
3: **for** $l = (L-1) : 1$ **do**
4: $\quad \delta_l := \delta_{l+1} W_{l+1}^\top \Theta(\mathbf{u}_l)$
5: $\quad \nabla_{w_{ijl}} \phi_W(\mathbf{x}) := \delta_l v_{i,l-1} \mathbf{e}_j \qquad \forall i, j$
6: **end for**

$-W_{l+1}^\top \leftarrow \Theta(x) \leftarrow u_L$

# Conclusion

How to train ANNs?

▶ Stochastic gradient descent

How to compute gradient?

▶ ANN is a nested function
▶ Thus we compute the gradient via the chain rule
▶ Lead to a recursive algorithm: backpropagation

# Outlook

Advanced training algorithms:

▶ Adagrad

▶ Adam

▶ Nesterov momentum