# Neural Attention Mechanisms

Guest Lecture: Deep Structured Prediction

Vlad Niculae

🐦 *@vnfrombucharest*

# Sequence-to-Sequence With Attention

```
words = [21, 79, 14]  # indices

embed = Embedding(vocab_sz, dim)
E = embed(words)  # (3 × dim)
enc = LSTM(dim, dim)
H = enc(E)  # (3 × dim)

dec = LSTM(2 * dim, dim)
initialize k=1, q[0], y[0]

while not done:
  s = H @ q[k - 1]  # attn scores
  # s = [-.3, -1.0, 1.8]

  p = softmax(s)  # attn proba
  # p = [.10, .05, .85]

  a = p @ H  # (1 × dim)

  q[k] = dec(a, y[k - 1], q[k - 1])
  y[k] = W_out @ q[k] + b
  k += 1
```

*United   Nations elections*

# Sequence-to-Sequence With Attention

```
words = [21, 79, 14]  # indices

embed = Embedding(vocab_sz, dim)
E = embed(words)  # (3 × dim)
enc = LSTM(dim, dim)
H = enc(E)  # (3 × dim)

dec = LSTM(2 * dim, dim)
initialize k=1, q[0], y[0]

while not done:
  s = H @ q[k - 1]  # attn scores
  # s = [-.3, -1.0, 1.8]

  p = softmax(s)  # attn proba
  # p = [.10, .05, .85]

  a = p @ H  # (1 × dim)

  q[k] = dec(a, y[k - 1], q[k - 1])
  y[k] = W_out @ q[k] + b
  k += 1
```

Encoder

$e_j$

*United  Nations elections*

# Sequence-to-Sequence With Attention

```
words = [21, 79, 14]  # indices

embed = Embedding(vocab_sz, dim)
E = embed(words)  # (3 × dim)
enc = LSTM(dim, dim)
H = enc(E)  # (3 × dim)

dec = LSTM(2 * dim, dim)
initialize k=1, q[0], y[0]

while not done:
  s = H @ q[k - 1]  # attn scores
  # s = [-.3, -1.0, 1.8]

  p = softmax(s)  # attn proba
  # p = [.10, .05, .85]

  a = p @ H  # (1 × dim)

  q[k] = dec(a, y[k - 1], q[k - 1])
  y[k] = W_out @ q[k] + b
  k += 1
```
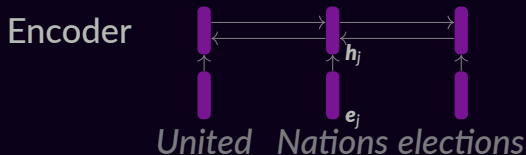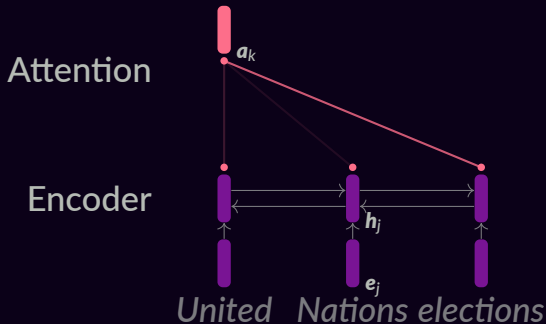
Encoder

$h_j$

$e_j$

*United Nations elections*

# Sequence-to-Sequence With Attention



Attention

Encoder

*United  Nations elections*

```
words = [21, 79, 14]  # indices

embed = Embedding(vocab_sz, dim)
E = embed(words)  # (3 × dim)
enc = LSTM(dim, dim)
H = enc(E)  # (3 × dim)

dec = LSTM(2 * dim, dim)
initialize k=1, q[0], y[0]

while not done:
  s = H @ q[k - 1]  # attn scores
  # s = [-.3, -1.0, 1.8]

  p = softmax(s)  # attn proba
  # p = [.10, .05, .85]

  a = p @ H  # (1 × dim)

  q[k] = dec(a, y[k - 1], q[k - 1])
  y[k] = W_out @ q[k] + b
  k += 1
```
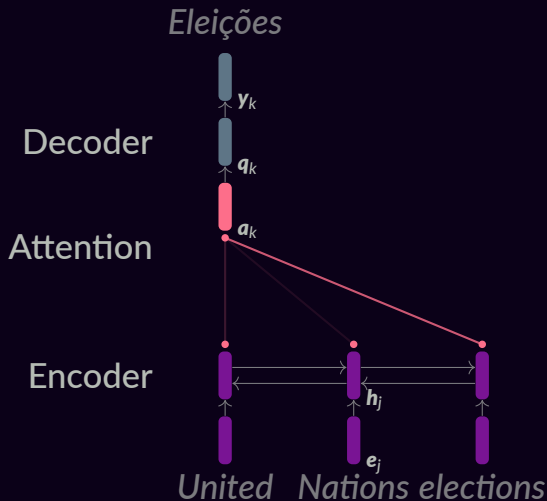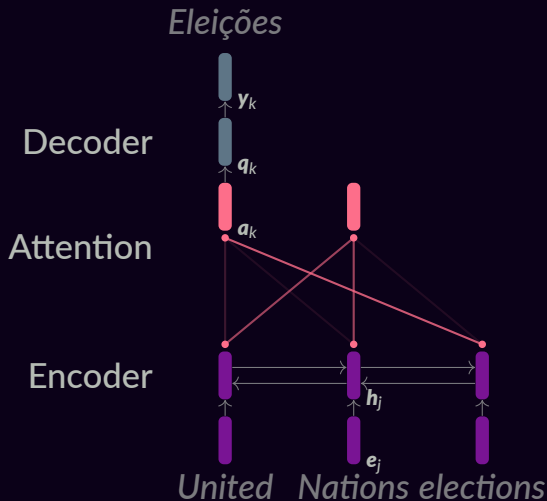
# Sequence-to-Sequence With Attention

(Bahdanau et al., 2015)



*Eleições*

Decoder $y_k$ $q_k$

Attention $a_k$

Encoder $h_j$ $e_j$

*United Nations elections*

```
words = [21, 79, 14]  # indices

embed = Embedding(vocab_sz, dim)
E = embed(words)  # (3 × dim)
enc = LSTM(dim, dim)
H = enc(E)  # (3 × dim)

dec = LSTM(2 * dim, dim)
initialize k=1, q[0], y[0]

while not done:
  s = H @ q[k - 1]  # attn scores
  # s = [-.3, -1.0, 1.8]

  p = softmax(s)  # attn proba
  # p = [.10, .05, .85]

  a = p @ H  # (1 × dim)

  q[k] = dec(a, y[k - 1], q[k - 1])
  y[k] = W_out @ q[k] + b
  k += 1
```
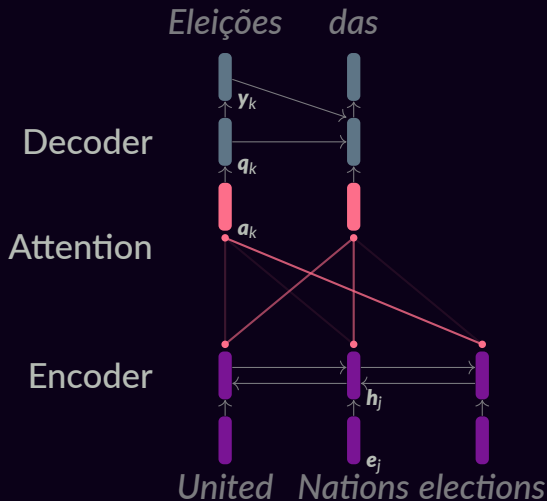
# Sequence-to-Sequence With Attention

(Bahdanau et al., 2015)

*Eleições*

Decoder

Attention

Encoder

*United  Nations elections*

```
words = [21, 79, 14]  # indices

embed = Embedding(vocab_sz, dim)
E = embed(words)  # (3 × dim)
enc = LSTM(dim, dim)
H = enc(E)  # (3 × dim)

dec = LSTM(2 * dim, dim)
initialize k=1, q[0], y[0]

while not done:
  s = H @ q[k - 1]  # attn scores
  # s = [-.3, -1.0, 1.8]

  p = softmax(s)  # attn proba
  # p = [.10, .05, .85]

  a = p @ H  # (1 × dim)

  q[k] = dec(a, y[k - 1], q[k - 1])
  y[k] = W_out @ q[k] + b
  k += 1
```
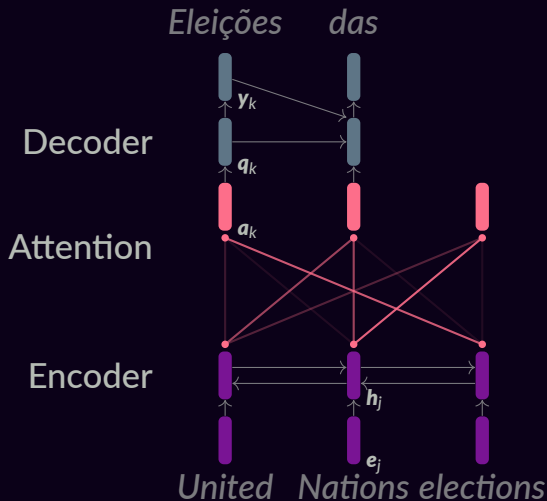
# Sequence-to-Sequence With Attention



```
words = [21, 79, 14]  # indices

embed = Embedding(vocab_sz, dim)
E = embed(words)  # (3 × dim)
enc = LSTM(dim, dim)
H = enc(E)  # (3 × dim)

dec = LSTM(2 * dim, dim)
initialize k=1, q[0], y[0]

while not done:
  s = H @ q[k - 1]  # attn scores
  # s = [-.3, -1.0, 1.8]

  p = softmax(s)  # attn proba
  # p = [.10, .05, .85]

  a = p @ H  # (1 × dim)

  q[k] = dec(a, y[k - 1], q[k - 1])
  y[k] = W_out @ q[k] + b
  k += 1
```

Figure labels: Decoder, Attention, Encoder

Eleições  das

$y_k$

$q_k$

$a_k$

$h_j$

$e_j$

United  Nations elections

# Sequence-to-Sequence With Attention

```
words = [21, 79, 14]  # indices

embed = Embedding(vocab_sz, dim)
E = embed(words)  # (3 × dim)
enc = LSTM(dim, dim)
H = enc(E)  # (3 × dim)

dec = LSTM(2 * dim, dim)
initialize k=1, q[0], y[0]

while not done:
  s = H @ q[k - 1]  # attn scores
  # s = [-.3, -1.0, 1.8]

  p = softmax(s)  # attn proba
  # p = [.10, .05, .85]

  a = p @ H  # (1 × dim)

  q[k] = dec(a, y[k - 1], q[k - 1])
  y[k] = W_out @ q[k] + b
  k += 1
```
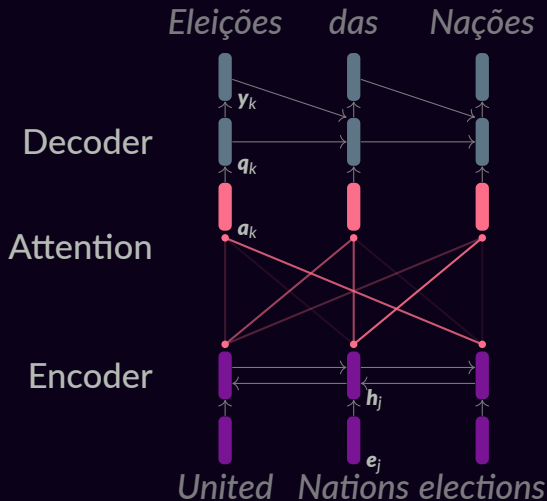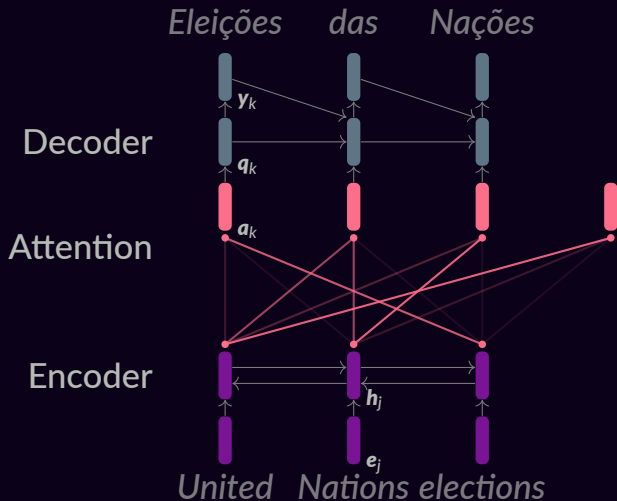
# Sequence-to-Sequence With Attention

*Eleições*  *das*  *Nações*

Decoder $y_k$ $q_k$

Attention $a_k$

Encoder $h_j$ $e_j$

*United  Nations elections*

```
words = [21, 79, 14]  # indices

embed = Embedding(vocab_sz, dim)
E = embed(words)  # (3 × dim)
enc = LSTM(dim, dim)
H = enc(E)  # (3 × dim)

dec = LSTM(2 * dim, dim)
initialize k=1, q[0], y[0]

while not done:
  s = H @ q[k - 1]  # attn scores
  # s = [-.3, -1.0, 1.8]

  p = softmax(s)  # attn proba
  # p = [.10, .05, .85]

  a = p @ H  # (1 × dim)

  q[k] = dec(a, y[k - 1], q[k - 1])
  y[k] = W_out @ q[k] + b
  k += 1
```
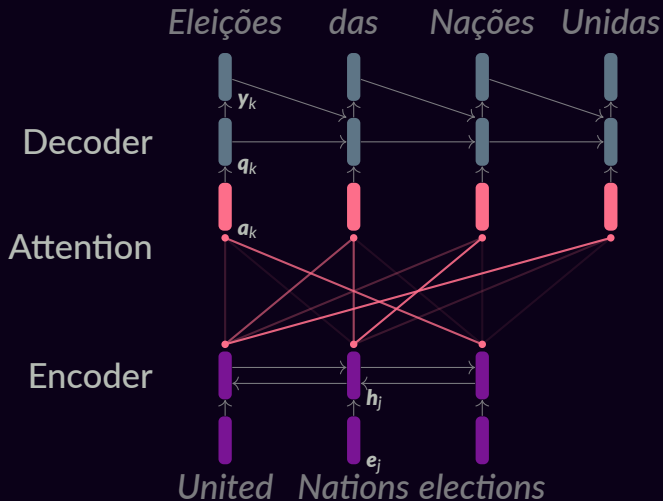
# Sequence-to-Sequence With Attention

```
words = [21, 79, 14]  # indices

embed = Embedding(vocab_sz, dim)
E = embed(words)  # (3 × dim)
enc = LSTM(dim, dim)
H = enc(E)  # (3 × dim)

dec = LSTM(2 * dim, dim)
initialize k=1, q[0], y[0]

while not done:
  s = H @ q[k - 1]  # attn scores
  # s = [-.3, -1.0, 1.8]

  p = softmax(s)  # attn proba
  # p = [.10, .05, .85]

  a = p @ H  # (1 × dim)

  q[k] = dec(a, y[k - 1], q[k - 1])
  y[k] = W_out @ q[k] + b
  k += 1
```
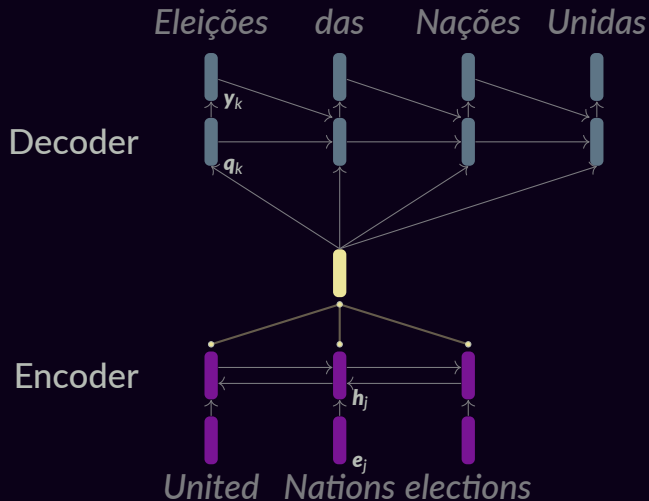
# Sequence-to-Sequence With Attention



```
words = [21, 79, 14]  # indices

embed = Embedding(vocab_sz, dim)
E = embed(words)  # (3 × dim)
enc = LSTM(dim, dim)
H = enc(E)  # (3 × dim)

dec = LSTM(2 * dim, dim)
initialize k=1, q[0], y[0]

while not done:
  s = H @ q[k - 1]  # attn scores
  # s = [-.3, -1.0, 1.8]

  p = softmax(s)  # attn proba
  # p = [.10, .05, .85]

  a = p @ H  # (1 × dim)

  q[k] = dec(a, y[k - 1], q[k - 1])
  y[k] = W_out @ q[k] + b
  k += 1
```
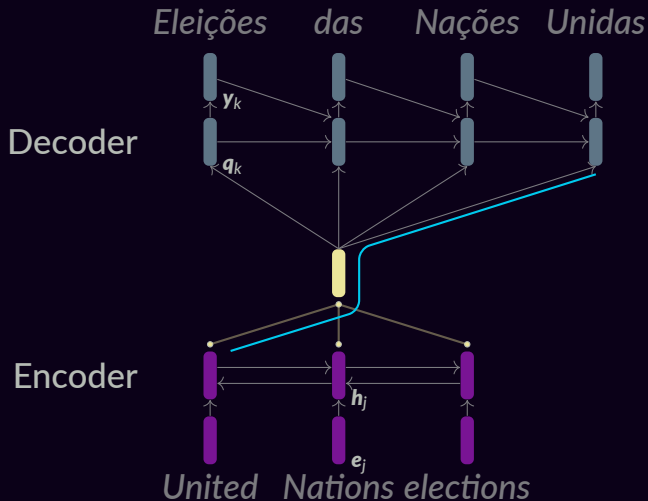
# Sequence-to-Sequence With Attention

(Bahdanau et al., 2015)

*Eleições  das  Nações  Unidas*

Decoder

$y_k$

$q_k$

Encoder

$h_j$

$e_j$

*United  Nations elections*

```
words = [21, 79, 14]  # indices

embed = Embedding(vocab_sz, dim)
E = embed(words)  # (3 × dim)
enc = LSTM(dim, dim)
H = enc(E)  # (3 × dim)

dec = LSTM(2 * dim, dim)
initialize k=1, q[0], y[0]

while not done:
  s = H @ q[k - 1]  # attn scores
  # s = [-.3, -1.0, 1.8]

  p = softmax(s)  # attn proba
  # p = [.10, .05, .85]

  a = p @ H  # (1 × dim)

  q[k] = dec(a, y[k - 1], q[k - 1])
  y[k] = W_out @ q[k] + b
  k += 1
```

# Sequence-to-Sequence With Attention

*Eleições*    *das*    *Nações*    *Unidas*

Decoder

$y_k$

$q_k$

Encoder

$h_j$

$e_j$

*United*   *Nations elections*

```
words = [21, 79, 14]  # indices

embed = Embedding(vocab_sz, dim)
E = embed(words)  # (3 × dim)
enc = LSTM(dim, dim)
H = enc(E)  # (3 × dim)

dec = LSTM(2 * dim, dim)
initialize k=1, q[0], y[0]

while not done:
  s = H @ q[k - 1]  # attn scores
  # s = [-.3, -1.0, 1.8]

  p = softmax(s)  # attn proba
  # p = [.10, .05, .85]

  a = p @ H  # (1 × dim)

  q[k] = dec(a, y[k - 1], q[k - 1])
  y[k] = W_out @ q[k] + b
  k += 1
```

# Attention as a shortcut

Attention doesn't make models more expressive,
it makes it easier to express "better" functions.

"You May Not Need Attention" for NMT,
but reordering is needed for good results.

(Press and Smith, 2018)

```python
# attention scores:
s = H @ W_attn @ state
# s = [-.3, -1.0, 1.8]

p = softmax(s)
# p = [.10, .05, .85]
```

*record scratch*

*freeze frame*

# 1. How to select an item from a set?

# How to select an item from a set?

*United*

*Nations*

*Elections*

# How to select an item from a set?

$$c_1$$
$$c_2$$
$$\dots$$
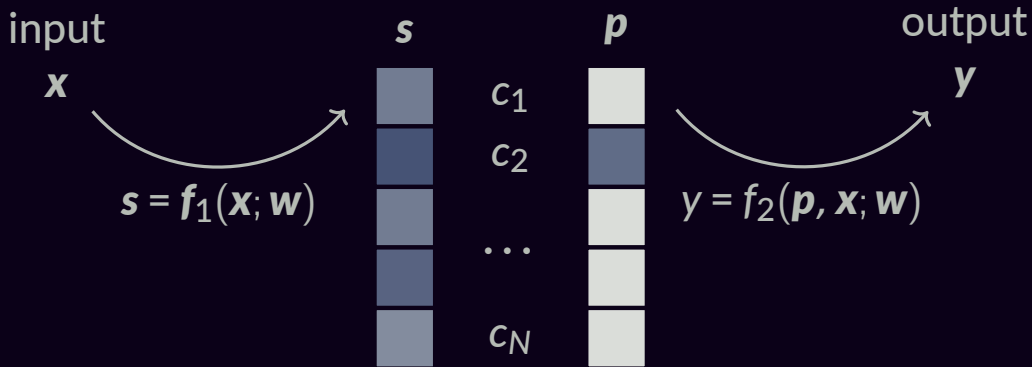$$c_N$$

# How to select an item from a set?

# How to select an item from a set?



$s$       $p$

$c_1$

$c_2$

$\ldots$

$c_N$

# How to select an item from a set?

input $\boldsymbol{x}$    $\boldsymbol{s}$    $\boldsymbol{p}$    output $\boldsymbol{y}$

$\boldsymbol{s} = \boldsymbol{f}_1(\boldsymbol{x}; \boldsymbol{w})$

$c_1$

$c_2$

$\dots$

$c_N$

$y = f_2(\boldsymbol{p}, \boldsymbol{x}; \boldsymbol{w})$

# How to select an item from a set?



input
$x$

$s$

$p$

output
$y$

$c_1$
$c_2$
$\ldots$
$c_N$

$s = f_1(x; w)$

$y = f_2(p, x; w)$

$\frac{\partial y}{\partial w} = ?$

# How to select an item from a set?

input $x$

$s = f_1(x; w)$

$s$

$c_1$
$c_2$
$\dots$
$c_N$

$p$

$y = f_2(p, x; w)$

output $y$

$\frac{\partial y}{\partial w} = ?$    or, essentially,    $\frac{\partial p}{\partial s} = ?$

# Winner Takes It All



$s$       $p$

$c_1$

$c_2$

...

$c_N$

$$\frac{\partial p}{\partial s} = ?$$

# Winner Takes It All



$s$     $p$

$c_1$

$c_2$

$\dots$

$c_N$

$$\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{s}} = ?$$

# Winner Takes It All



$s$       $p$

$c_1$

$c_2$

$\ldots$

$c_N$

$$\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{s}} = ?$$

# Winner Takes It All



$s$       $p$

$c_1$

$c_2$

$\ldots$

$c_N$

$$\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{s}} = ?$$

# Winner Takes It All

# Winner Takes It All

# Winner Takes It All

# Winner Takes It All

# Argmax



$$\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{s}} = 0$$

# Argmax vs. Softmax

$$p_j = \exp(s_j)/Z$$



$$s \qquad p$$

$c_1$

$c_2$

$\cdots$

$c_N$

$$\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{s}} = \mathrm{diag}(\boldsymbol{p}) - \boldsymbol{p}\boldsymbol{p}^\top$$

# Background: Optimization

$$f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$$

$$\min_{\boldsymbol{x}} f(\boldsymbol{x}) := v \text{ s.t. } (a) \; \exists \boldsymbol{x}^\star \in \mathbb{R}^d, f(\boldsymbol{x}^\star) = v$$

$$(b) \; \forall \boldsymbol{x}' \in \mathbb{R}^d, f(\boldsymbol{x}') \geq v$$

$$\arg\min_{\boldsymbol{x}} f(\boldsymbol{x}) := \{\boldsymbol{x}^\star \in \mathbb{R}^d : f(\boldsymbol{x}^\star) = \min_{\boldsymbol{x}} f(\boldsymbol{x})\}$$

$f$ convex: optimization algos available

$f$ **strictly** convex: $\arg\min_{\boldsymbol{x}} f(\boldsymbol{x}) = \{\boldsymbol{x}^\star\}$

# Background: Constrained Optimization

$$\min_{\boldsymbol{x} \in \mathcal{X} \subset \mathbb{R}^d} f(\boldsymbol{x})$$

The indicator function: $\mathrm{Id}_{\mathcal{X}}(\boldsymbol{x}) = \begin{cases} 0, & \boldsymbol{x} \in \mathcal{X}, \\ \infty, & \boldsymbol{x} \notin \mathcal{X}. \end{cases}$

$$\arg\min_{\boldsymbol{x} \in \mathcal{X}} f(\boldsymbol{x}) = \arg\min_{\boldsymbol{x} \in \mathbb{R}^d} f(\boldsymbol{x}) + \mathrm{Id}_{\mathcal{X}}(\boldsymbol{x}).$$
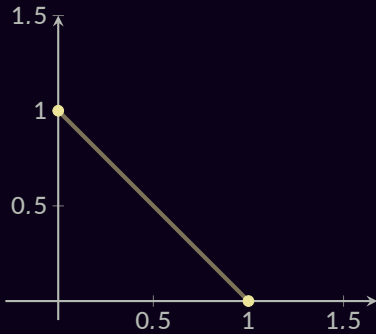
$\mathrm{Id}_{\mathcal{X}}$ is a convex function when $\mathcal{X}$ a convex set.

# The Simplex

$$\triangle = \{ \boldsymbol{p} \in \mathbb{R}^d : \ \boldsymbol{p} \geq \boldsymbol{0}, \ \boldsymbol{1}^\top \boldsymbol{p} = 1 \}$$

# The Simplex

$$\triangle = \{ \boldsymbol{p} \in \mathbb{R}^d : \boldsymbol{p} \geq \boldsymbol{0}, \ \boldsymbol{1}^\top \boldsymbol{p} = 1 \}$$



$d = 2$

# The Simplex

$$\triangle = \{\boldsymbol{p} \in \mathbb{R}^d : \boldsymbol{p} \geq \boldsymbol{0}, \ \boldsymbol{1}^\top \boldsymbol{p} = 1\}$$



$\boldsymbol{p} = [1, 0]$

$d = 2$

# The Simplex

$$\triangle = \{ \boldsymbol{p} \in \mathbb{R}^d : \boldsymbol{p} \geq \boldsymbol{0}, \ \boldsymbol{1}^\top \boldsymbol{p} = 1 \}$$



$d = 2$

# The Simplex

$$\triangle = \{ \boldsymbol{p} \in \mathbb{R}^d : \boldsymbol{p} \geq \boldsymbol{0}, \ \boldsymbol{1}^\top \boldsymbol{p} = 1 \}$$



$\boldsymbol{p} = [0, 1]$

$\boldsymbol{p} = [\frac{1}{2}, \frac{1}{2}]$

$\boldsymbol{p} = [1, 0]$
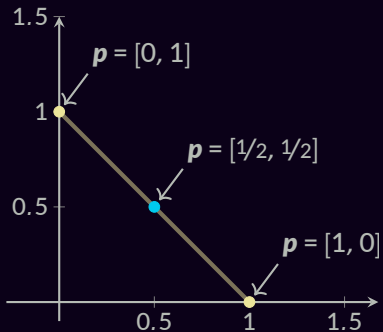
$d = 2$

# The Simplex

$$\triangle = \{ p \in \mathbb{R}^d : p \geq 0, \ 1^\top p = 1 \}$$



$d = 2$

$n = 3$

# The Simplex

$$\triangle = \{ \boldsymbol{p} \in \mathbb{R}^d : \boldsymbol{p} \geq \boldsymbol{0}, \ \boldsymbol{1}^\top \boldsymbol{p} = 1 \}$$



$d = 2$

$n = 3$

# The Simplex

$$\triangle = \{ \boldsymbol{p} \in \mathbb{R}^d : \boldsymbol{p} \geq \boldsymbol{0}, \ \boldsymbol{1}^\top \boldsymbol{p} = 1 \}$$



$d = 2$

$n = 3$

# The Simplex

$$\triangle = \{ p \in \mathbb{R}^d : p \geq 0, \ 1^\top p = 1 \}$$



$p = [0, 1]$

$p = [1/2, 1/2]$

$p = [1, 0]$

$d = 2$

$p = [0, 0, 1]$

$p = [1/3, 1/3, 1/3]$

$p = [0, 1, 0]$

$n = 3$

# Highest Element of a Vector

$$\max_{j} s_j = \max_{p \in \triangle} p^\top s$$

Fundamental Thm. Lin. Prog.
(Dantzig et al., 1955)



$d = 2$

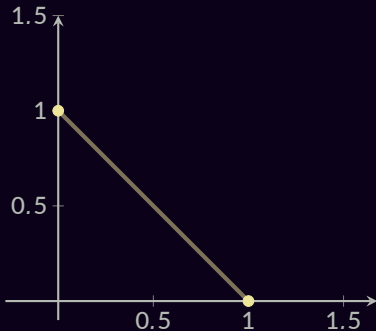$n = 3$

# Highest Element of a Vector

$$\max_{j} s_j = \max_{p \in \triangle} p^\top s$$

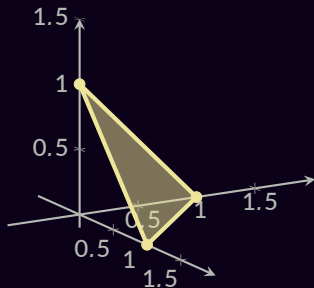Fundamental Thm. Lin. Prog.
(Dantzig et al., 1955)



$s = [.2, 1.4]$

$d = 2$

$n = 3$

# Highest Element of a Vector

$$\max_j s_j = \max_{p \in \triangle} p^\top s$$

Fundamental Thm. Lin. Prog.
(Dantzig et al., 1955)
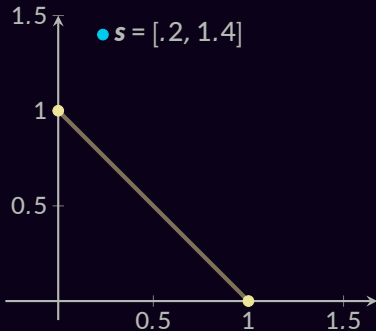


$s = [.2, 1.4]$

$d = 2$

$n = 3$

# Highest Element of a Vector

$$\max_j s_j = \max_{p \in \triangle} p^\top s$$

Fundamental Thm. Lin. Prog.
(Dantzig et al., 1955)



$d = 2$

$n = 3$

# Highest Element of a Vector

$$\max_j s_j = \max_{p \in \triangle} p^\top s$$

Fundamental Thm. Lin. Prog.
(Dantzig et al., 1955)



$s = [.2, 1.4]$

$p^\star = [0, 1]$
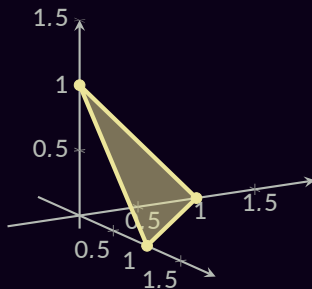
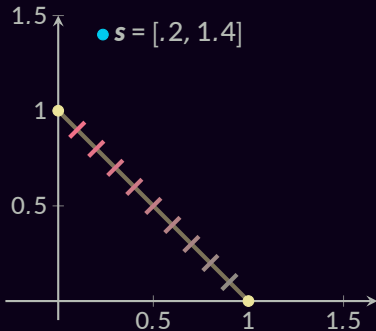$d = 2$

$s = [.7, .1, 1.5]$

$n = 3$

# Highest Element of a Vector

$$\max_j s_j = \max_{p \in \triangle} p^\top s$$

Fundamental Thm. Lin. Prog.
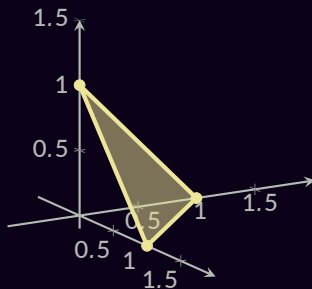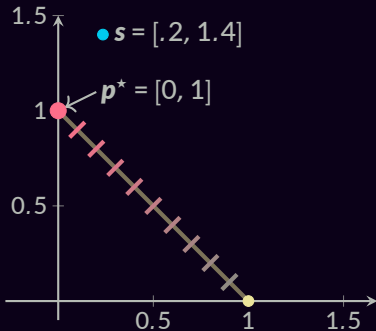(Dantzig et al., 1955)



$d = 2$

$n = 3$

# Highest Element of a Vector

$$\max_{j} s_j = \max_{\boldsymbol{p} \in \triangle} \boldsymbol{p}^\top \boldsymbol{s}$$

Fundamental Thm. Lin. Prog.
(Dantzig et al., 1955)



$\boldsymbol{s} = [.2, 1.4]$

$\boldsymbol{p}^\star = [0, 1]$

$d = 2$



$\boldsymbol{s} = [.7, .1, 1.5]$

$\boldsymbol{p}^\star = [0, 0, 1]$

$n = 3$

# Smoothed Max Operators

$$\max_{\Omega}(\boldsymbol{s}) = \max_{\boldsymbol{p}\in\triangle} \boldsymbol{p}^{\top}\boldsymbol{s} - \Omega(\boldsymbol{p})$$

# Smoothed Max Operators

$$\max{}_{\Omega}(\boldsymbol{s}) = \max_{\boldsymbol{p}\in\triangle} \boldsymbol{p}^{\top}\boldsymbol{s} - \Omega(\boldsymbol{p})$$

- argmax: $\Omega(\boldsymbol{p}) = 0$

# Smoothed Max Operators

$$\max{}_\Omega(\boldsymbol{s}) = \max_{\boldsymbol{p} \in \triangle} \boldsymbol{p}^\top \boldsymbol{s} - \Omega(\boldsymbol{p})$$

- argmax: $\Omega(\boldsymbol{p}) = 0$

- softmax: $\Omega(\boldsymbol{p}) = \sum_j p_j \log p_j$

# Smoothed Max Operators

$$\max\nolimits_{\Omega}(\boldsymbol{s}) = \max_{\boldsymbol{p} \in \triangle} \boldsymbol{p}^{\top}\boldsymbol{s} - \Omega(\boldsymbol{p})$$



- ●    argmax:   $\Omega(\boldsymbol{p}) = 0$

- ●    softmax:   $\Omega(\boldsymbol{p}) = \sum_j p_j \log p_j$

- ●    sparsemax:   $\Omega(\boldsymbol{p}) = \tfrac{1}{2}\|\boldsymbol{p}\|_2^2$

(Martins and Astudillo, 2016)

softmax

sparsemax

# Sparsemax

$$\text{sparsemax}(\boldsymbol{s}) = \underset{\boldsymbol{p} \in \triangle}{\arg\max}\, \boldsymbol{p}^\top \boldsymbol{s} - \tfrac{1}{2}\|\boldsymbol{p}\|_2^2$$

$$= \underset{\boldsymbol{p} \in \triangle}{\arg\min}\, \|\boldsymbol{p} - \boldsymbol{s}\|_2^2$$

**Computation:**

$$\boldsymbol{p}^\star = [\boldsymbol{s} - \tau\boldsymbol{1}]_+$$

$$s_i > s_j \Rightarrow p_i \geq p_j$$

$O(d)$ via partial sort

(Held et al., 1974; Brucker, 1984; Condat, 2016)

**Backward pass:**

$$\boldsymbol{J}_{\text{sparsemax}} = \text{diag}(\boldsymbol{s}) - \tfrac{1}{|S|}\boldsymbol{s}\boldsymbol{s}^\top$$

$$\text{where } \mathcal{S} = \{j : p_j^\star > 0\},$$

$$s_j = [\![ j \in \mathcal{S} ]\!]$$

(Martins and Astudillo, 2016)

fusedmax

# Smoothed Max Operators

$$\max_{\Omega}(\boldsymbol{s}) = \max_{\boldsymbol{p} \in \triangle} \boldsymbol{p}^{\top} \boldsymbol{s} - \Omega(\boldsymbol{p})$$

- ● argmax: $\Omega(\boldsymbol{p}) = 0$

- ● softmax: $\Omega(\boldsymbol{p}) = \sum_j p_j \log p_j$

- ● sparsemax: $\Omega(\boldsymbol{p}) = \frac{1}{2}\|\boldsymbol{p}\|_2^2$

# Fusedmax

$$\text{fusedmax}(\boldsymbol{s}) = \underset{\boldsymbol{p} \in \triangle}{\arg\max} \, \boldsymbol{p}^\top \boldsymbol{s} - \tfrac{1}{2}\|\boldsymbol{p}\|_2^2 - \sum_{2 \le j \le d} |p_j - p_{j-1}|$$

$$= \underset{\boldsymbol{p} \in \triangle}{\arg\min} \, \|\boldsymbol{p} - \boldsymbol{s}\|_2^2 + \sum_{2 \le j \le d} |p_j - p_{j-1}|$$

$$\text{prox}_{\text{fused}}(\boldsymbol{s}) = \underset{\boldsymbol{p} \in \mathbb{R}^d}{\arg\min} \, \|\boldsymbol{p} - \boldsymbol{s}\|_2^2 + \sum_{2 \le j \le d} |p_j - p_{j-1}|$$

**Proposition:** $\text{fusedmax}(\boldsymbol{s}) = \text{sparsemax}\big(\text{prox}_{\text{fused}}(\boldsymbol{s})\big)$

(Niculae and Blondel, 2017)

fusedmax $p_j - p_{j-1}|$

-1|

prox$_\text{fused}$ -1|

**Propos** used($s$))



"Fused Lasso" a.k.a. 1-d Total Variation

(Tibshirani et al., 2005)

fusedmax

# Constrained Attention

$$\arg\max_{\boldsymbol{p}\in\triangle\cup\mathcal{X}} \boldsymbol{p}^\top\boldsymbol{s} - \Omega(\boldsymbol{p})$$

$$= \arg\max_{\boldsymbol{p}\in\triangle} \boldsymbol{p}^\top\boldsymbol{s} - \underbrace{\Omega_\mathcal{X}(\boldsymbol{p})}_{\Omega+\mathsf{Id}_\mathcal{X}}$$

# Constrained Attention

$$\arg\max_{\boldsymbol{p}\in\triangle\cup\mathcal{X}} \boldsymbol{p}^\top\boldsymbol{s} - \Omega(\boldsymbol{p})$$

$$= \arg\max_{\boldsymbol{p}\in\triangle} \boldsymbol{p}^\top\boldsymbol{s} - \underbrace{\Omega_\mathcal{X}(\boldsymbol{p})}_{\Omega + \mathsf{Id}_\mathcal{X}}$$

**Example:** upper bounds $\mathcal{X} = \{\boldsymbol{p}\in\mathbb{R}^d : p_j \leq b_j\}$
constrained softmax (Martins and Kreutzer, 2017) and sparsemax (Malaviya et al., 2018)
Application: incorporating fertility in Neural MT

# Example: Source Sentence with Three Words

# Smoothed Max Operators

$$\max{}_{\Omega}(\boldsymbol{s}) = \max_{\boldsymbol{p} \in \triangle} \boldsymbol{p}^{\top} \boldsymbol{s} - \Omega(\boldsymbol{p})$$



- ●   argmax:   $\Omega(\boldsymbol{p}) = 0$

- ●   softmax:   $\Omega(\boldsymbol{p}) = \sum_j p_j \log p_j$

- ●   sparsemax:   $\Omega(\boldsymbol{p}) = \frac{1}{2}\|\boldsymbol{p}\|_2^2$

    fusedmax:   $\Omega(\boldsymbol{p}) = \frac{1}{2}\|\boldsymbol{p}\|_2^2 + \sum_j |p_j - p_{j-1}|$

    csparsemax:   $\Omega(\boldsymbol{p}) = \frac{1}{2}\|\boldsymbol{p}\|_2^2 + \mathsf{Id}_{\boldsymbol{p} \leq \boldsymbol{b}}$

# Smoothed Max Operators

$$\max\nolimits_{\Omega}(\boldsymbol{s}) = \max_{\boldsymbol{p} \in \triangle} \boldsymbol{p}^{\top}\boldsymbol{s} - \Omega(\boldsymbol{p})$$



- argmax: $\Omega(\boldsymbol{p}) = 0$

- softma

- sparsema

  fusedmax: $\Omega(\boldsymbol{p}) - \tfrac{1}{2}\|\boldsymbol{p}\|_2^- + \sum_j |p_j - p_{j-1}|$

  csparsemax: $\Omega(\boldsymbol{p}) = \tfrac{1}{2}\|\boldsymbol{p}\|_2^2 + \mathrm{Id}_{\boldsymbol{p} \leq \boldsymbol{b}}$

Black-box solvers available (e.g. FISTA), specialized solvers can be much faster.

$[.3, .2, .5]$

# 2. Attention architectures.

# Computing the scores

$$s_j = \sigma(\boldsymbol{h}_j, \boldsymbol{q})$$

| name | $\sigma(\boldsymbol{h}, \boldsymbol{q})$ | |
|---|---|---|
| additive | $\boldsymbol{v}^\top \tanh(\boldsymbol{W}_1\boldsymbol{h} + \boldsymbol{W}_2\boldsymbol{q})$ | (Bahdanau et al., 2015) |
| dot-product | $\boldsymbol{h}^\top \boldsymbol{q}$ | (Luong et al., 2015) |
| bilinear | $\boldsymbol{h}^\top \boldsymbol{W}\boldsymbol{q}$ | (Luong et al., 2015) |
| scaled | $(1/\sqrt{d})\, \boldsymbol{h}^\top \boldsymbol{W}\boldsymbol{q}$ | (Vaswani et al., 2017) |

# Beyond seq2seq

The spirit of *attention mechanisms* reaches far:

- ▶ Key-Value Attention
- ▶ Multi-head Attention
- ▶ Self-Attention and the Transformer
- ▶ Hierarchical Attention
- ▶ Memory Networks, Pointer Networks, Neural Turing Machines...

# Key-Value Attention

idea: the objects we average *(values)*
and the objects used to compute scores *(keys)*
don't need to be identical!

$$s_j = \boldsymbol{h}_j^\top \boldsymbol{q}$$

$$\boldsymbol{u} = \text{softmax}(\boldsymbol{s})^\top \boldsymbol{H}$$

$$s_j = \boldsymbol{k}_j^\top \boldsymbol{q}$$

$$\boldsymbol{u} = \text{softmax}(\boldsymbol{s})^\top \boldsymbol{V}$$

# Multi-head Attention

idea: compute $k$ different attention averages,
& concatenate the outputs.

$$s_j = \boldsymbol{k}_j^\top \boldsymbol{q}$$

$$\boldsymbol{u} = \text{softmax}(\boldsymbol{s})^\top \boldsymbol{V}$$

$$s_j^{(i)} = (\boldsymbol{W}_k^{(i)} \boldsymbol{k}_j)^\top (\boldsymbol{W}_q^{(i)} \boldsymbol{q})$$

$$\boldsymbol{u}^{(i)} = \text{softmax}(\boldsymbol{s}^{(i)})^\top (\boldsymbol{V} \boldsymbol{W}_v^{(i)})$$

$$\boldsymbol{u} = [\boldsymbol{u}^{(1)}; \cdots ; \boldsymbol{u}^{(k)}]$$

```
u = softmax(K @ q) @ V
```

```
for i in range(num_heads):
    Ki = K @ Wk[i].t()
    Vi = V @ Wv[i].t()
    qi = q @ Wq[i].t()
    ui = softmax(Ki @ qi) @ Vi
u = concat(ui)
```

# Self-attention

Attention as an *encoder layer*

...

*The   bears   eat   the   pretty   ones*

# Self-attention

Attention as an *encoder layer*

(Cheng et al., 2016)

…

**Transformer** (Vaswani et al., 2017): very deep self-attention
replacing LSTMs in encoder & decoder

*The   bears   eat   the   pretty   ones*

# Hierarchical Attention

Encode document by first encoding its sentences.

*s*

# References I

Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio (2015). "Neural machine translation by jointly learning to align and translate". In: *Proc. of ICLR*.
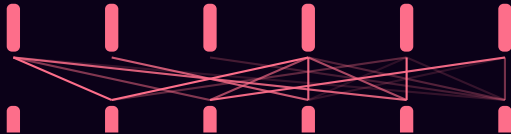
Brucker, Peter (1984). "An $O(n)$ algorithm for quadratic knapsack problems". In: *Operations Research Letters* 3.3, pp. 163–166.

Cheng, Jianpeng, Li Dong, and Mirella Lapata (2016). "Long Short-Term Memory-Networks for Machine Reading". In: *Proc. of EMNLP*.

Condat, Laurent (2016). "Fast projection onto the simplex and the $\ell_1$ ball". In: *Mathematical Programming* 158.1-2, pp. 575–585.

Dantzig, George B, Alex Orden, and Philip Wolfe (1955). "The generalized simplex method for minimizing a linear form under linear inequality restraints". In: *Pacific Journal of Mathematics* 5.2, pp. 183–195.

Graves, Alex, Greg Wayne, and Ivo Danihelka (2014). "Neural Turing Machines". In: *arXiv preprint arXiv:1410.5401*.

Held, Michael, Philip Wolfe, and Harlan P Crowder (1974). "Validation of subgradient optimization". In: *Mathematical Programming* 6.1, pp. 62–88.

Luong, Minh-Thang, Hieu Pham, and Christopher D Manning (2015). "Effective approaches to attention-based neural machine translation". In: *Proc. of EMNLP*.

Malaviya, Chaitanya, Pedro Ferreira, and André F. T. Martins (2018). "Sparse and constrained attention for neural machine translation". In: *Proc. of ACL*.

Martins, André FT and Ramón Fernandez Astudillo (2016). "From softmax to sparsemax: A sparse model of attention and multi-label classification". In: *Proc. of ICML*.

# References II

Martins, André FT and Julia Kreutzer (2017). "Learning What's Easy: Fully Differentiable Neural Easy-First Taggers". In: *Proc. of EMNLP*, pp. 349–362.

Niculae, Vlad and Mathieu Blondel (2017). "A regularized framework for sparse and structured neural attention". In: *Proc. of NeurIPS*.

Press, Ofir and Noah A Smith (2018). "You May Not Need Attention". In: *arXiv preprint arXiv:1810.13409*.

Sukhbaatar, Sainbayar, Jason Weston, and Rob Fergus (2015). "End-to-end memory networks". In: *Proc. of NeurIPS*.

Tibshirani, Robert, Michael Saunders, Saharon Rosset, Ji Zhu, and Keith Knight (2005). "Sparsity and smoothness via the fused lasso". In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 67.1, pp. 91–108.

Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin (2017). "Attention Is All You Need". In: *Proc. of NeurIPS*.

Yang, Zichao, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy (2016). "Hierarchical attention networks for document classification". In: *Proc. of NAACL-HLT*.