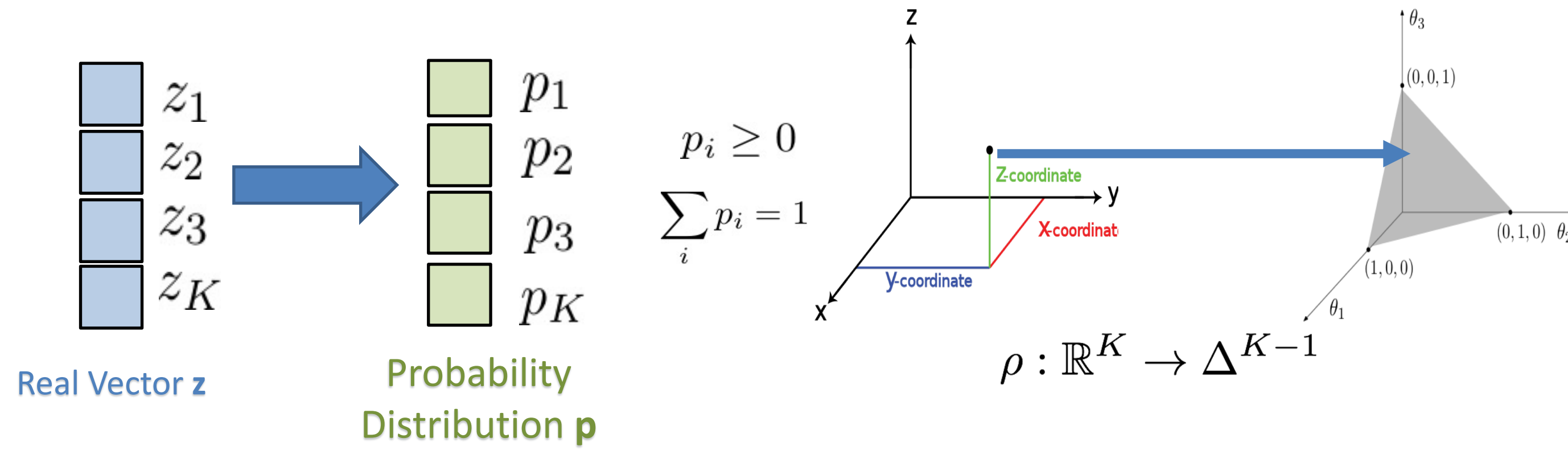


Probability Mapping Functions

- We are looking for a function ρ which takes a real vector \mathbf{z} and produces a probability distribution \mathbf{p} .



Applications:

1. **Probabilistic classification** – Multiclass classification, Multilabel Classification.
2. **Neural Attention Models** – Attention networks need a probability distribution over input states while generating output states.
3. **Memory Networks, Reinforcement Learning, Knowledge distillation** and many more.

Known probability mapping functions:

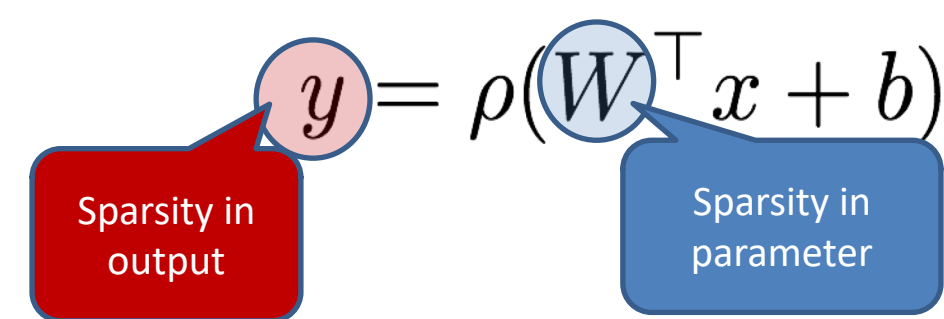
Limitations:

1. Softmax $\rho_i(\mathbf{z}) = \frac{\exp(z_i)}{\sum_{j \in [K]} \exp(z_j)}$ **1. Cannot be sparse**
2. Sum-normalization $\rho_i(\mathbf{z}) = \frac{z_i}{\sum_{j \in [K]} z_j}$ **2. Not full-domain**
3. Spherical softmax $\rho_i(\mathbf{z}) = \frac{z_i^2}{\sum_{j \in [K]} z_j^2}$ **3. Not monotonic**

Need for sparsity

1. In multilabel classification, ONLY A FEW labels out of 1000s of possible labels are TRUE.
2. In attention models/memory networks, sparser probabilities make COMPUTATION FASTER.

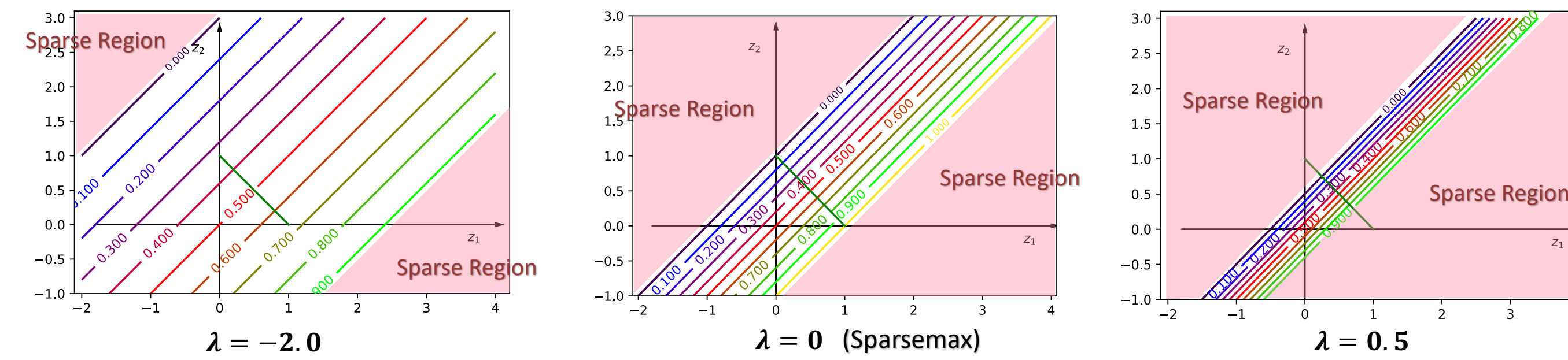
Sparsity in output VERSUS sparsity in model parameters.



Controls for Sparsity

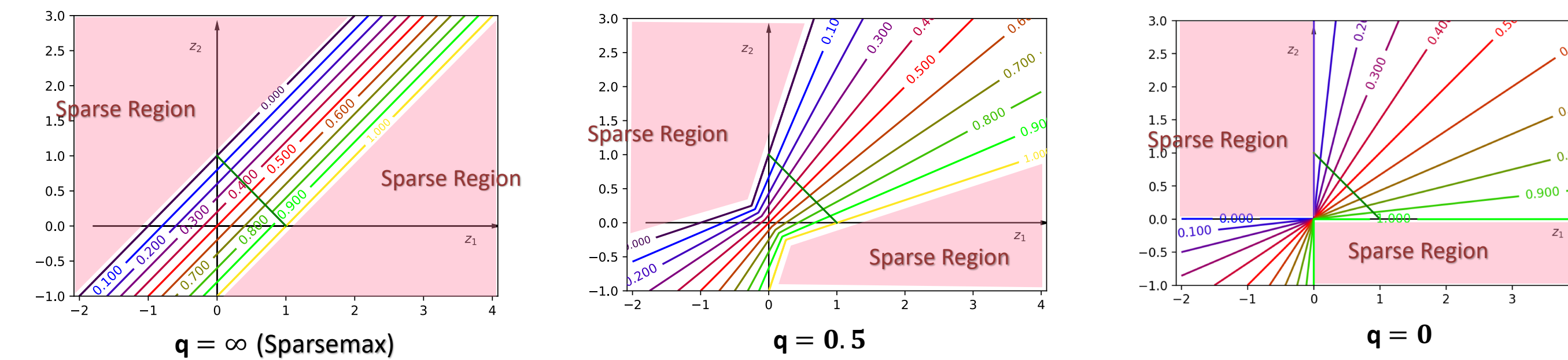
- Sparsegen-lin** (Control over width of non-sparse region):

$$\rho(\mathbf{z}) = \text{sparsegen-lin}(\mathbf{z}) = \underset{\mathbf{p} \in \Delta^{K-1}}{\operatorname{argmin}} \|\mathbf{p} - \mathbf{z}\|_2^2 - \lambda \|\mathbf{p}\|_2^2$$



- Sparse-hourglass** (Control over shape of non-sparse region):

$$\rho(\mathbf{z}) = \text{sparsehourglass}(\mathbf{z}) = \underset{\mathbf{p} \in \Delta^{K-1}}{\operatorname{argmin}} \left\| \mathbf{p} - \frac{1 + Kq}{|\sum_{i \in [K]} z_i| + Kq} \mathbf{z} \right\|_2^2$$



Translation Invariance!!!

Scale Invariance!!!

- Parameter q helps trade-off between translation and scale invariances.

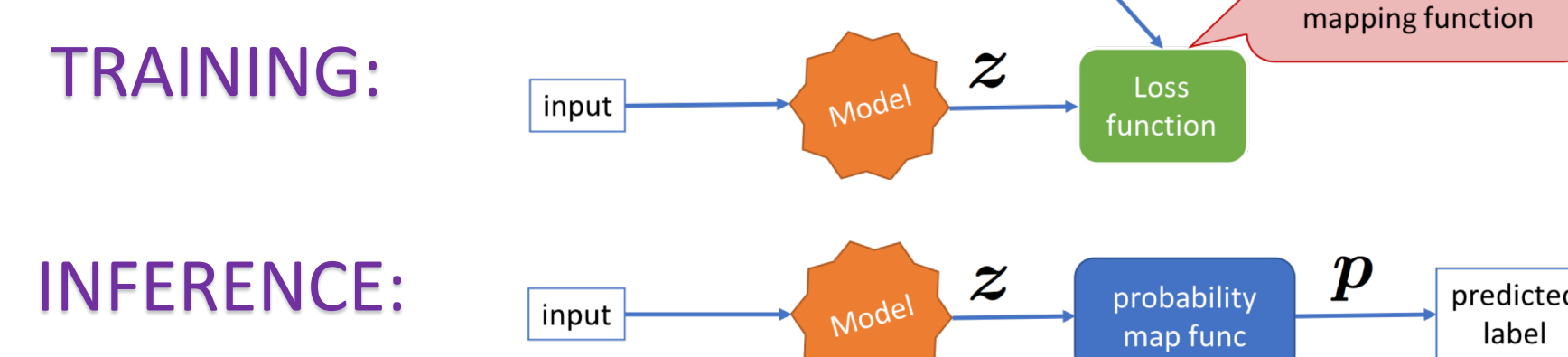
1. **Translation Invariance:** Adding a constant value to all dimensions of \mathbf{z} keeps \mathbf{p} unchanged.
2. **Scale Invariance:** Multiplying all dimensions of \mathbf{z} by a constant value keeps \mathbf{p} unchanged.

Sparsity Inducing Loss Functions

- Setting: Multilabel Classification**

- More than one labels for an instance can be true.
- *Usual approach* : Separate logistic sigmoid based binary classifier for every label followed by thresholding.
- *In this work*: Apply sparse probability mapping function. Non-zeroes are predicted labels, zeroes are non-labels.

- Training and Inference:



- Convex hinge-based loss functions:

$$\mathcal{L}_{\text{sparsehg,hinge}}(\mathbf{z}, \eta) = \sum_{i \neq 0, \eta_i \neq 0} |z_i - z_j| + \sum_{i \neq 0, \eta_j = 0} \max \left\{ \frac{\eta_i}{\hat{\alpha}(\mathbf{z})} - (z_i - z_j), 0 \right\} \quad \hat{\alpha}(\mathbf{z}) = \frac{1 + Kq}{|\sum_{i \in [K]} z_i| + Kq}$$

$$\mathcal{L}_{\text{sparsegen-lin,hinge}}(\mathbf{z}, \eta) = \frac{1}{1 - \lambda} \sum_{i \neq 0, \eta_i \neq 0} |z_i - z_j| + \sum_{i \neq 0, \eta_j = 0} \max \left\{ \eta_i - \frac{z_i - z_j}{1 - \lambda}, 0 \right\}$$

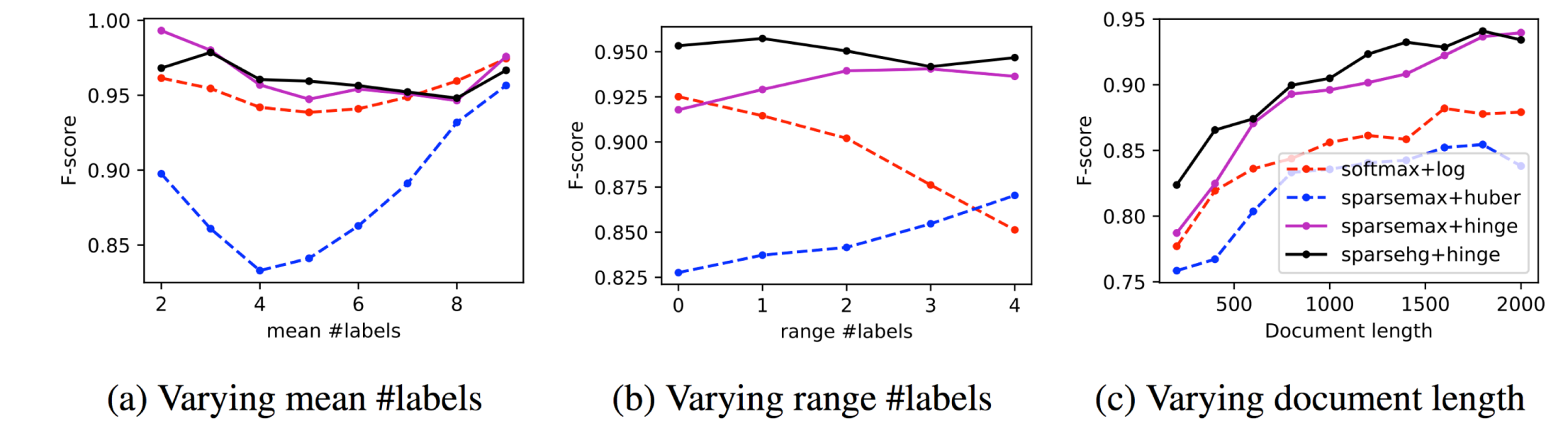
Multilabel Classification

Synthetic Multilabel Experimental Setup:

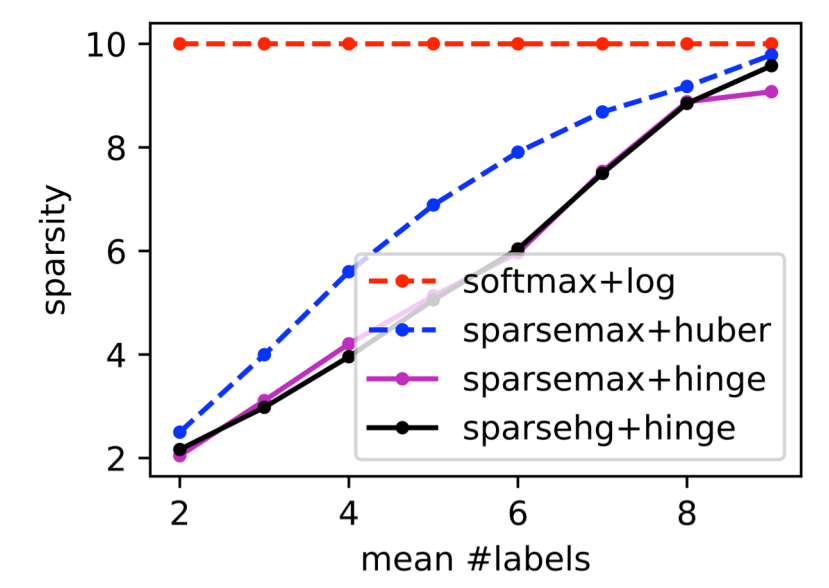
- a) Varying mean #labels.
- b) Varying range #labels.
- c) Vrying document length.

Competing models:

1. Baseline **softmax+log**.
2. Baseline **sparsemax+huber** [ICML 2016].
3. Proposed **sparsemax+hinge**.
4. Proposed **sparsehg+hinge**.



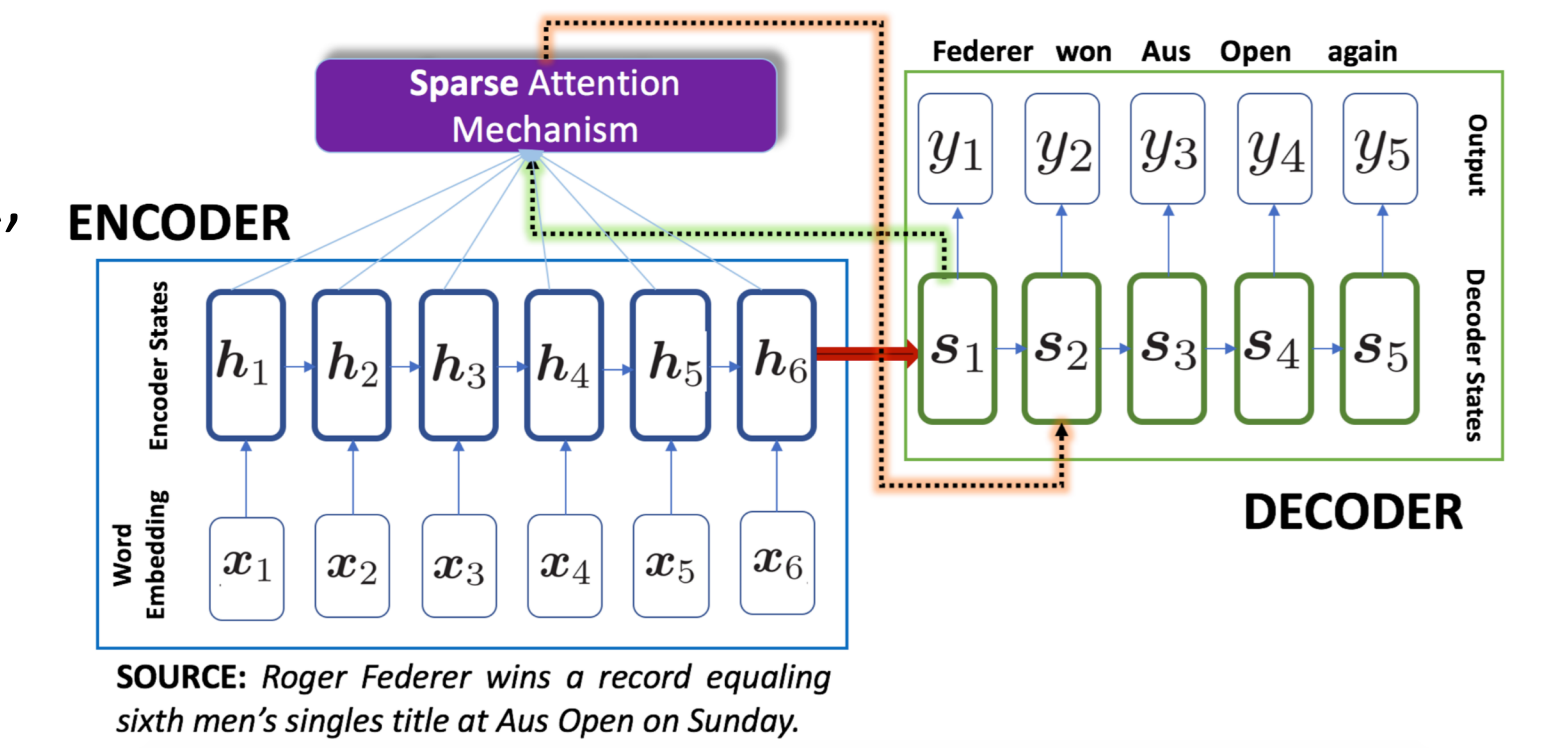
More accurate predictions
Sparser outputs



Controlled Sparse Attention

Seq2seq Models with attention:

- a) Neural Machine Translation (EN-FR, FR-EN).
- b) Abstractive Summarization (Gigaword, DUC2003, DUC2004).

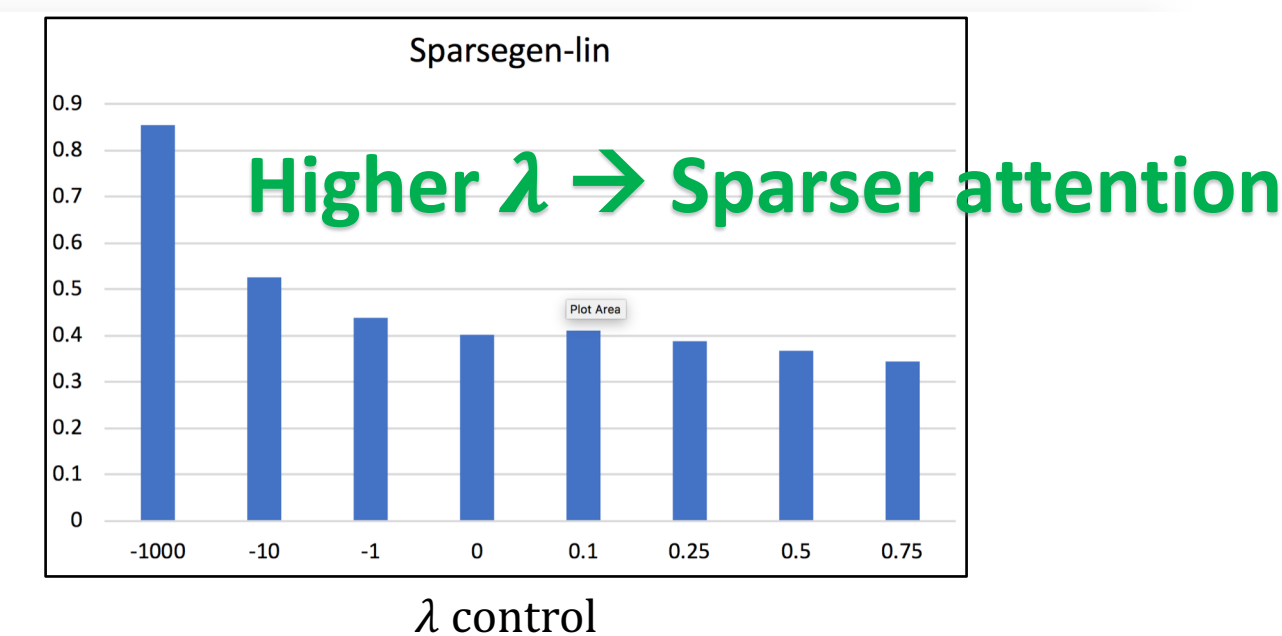


SOURCE: Roger Federer wins a record equalling sixth men's singles title at Aus Open on Sunday.

OpenNMT framework (PyTorch).

- Replace 'softmax' with **sparsemax**, **Sparsegen-lin** and **Sparsehourglass**.
- Also varied temperature in softmax as another baseline.

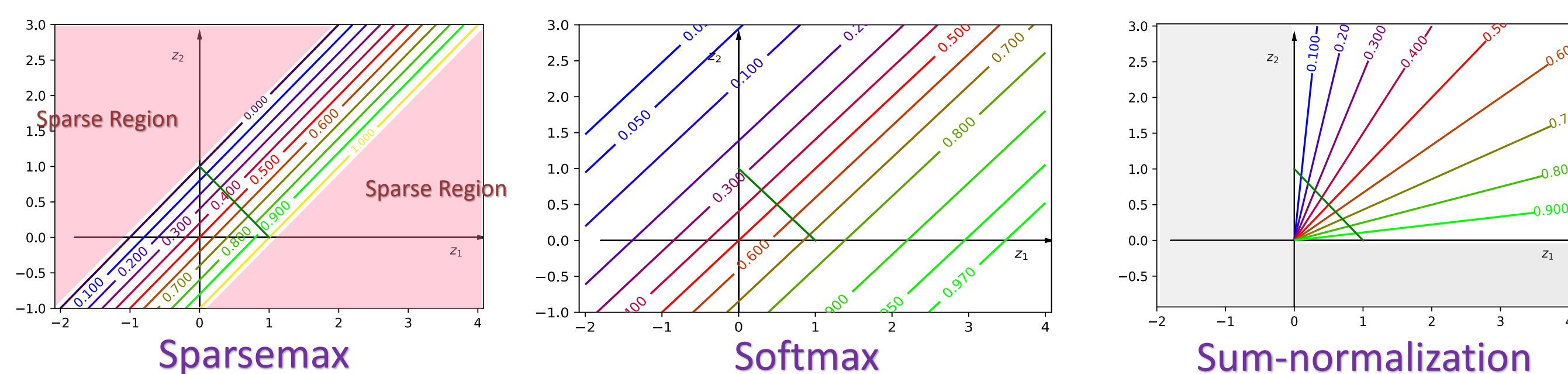
Summary of Results:



Sparse Probability Maps

- Sparsemax** – Projection onto simplex [ICML 2016]:

$$\rho(\mathbf{z}) = \underset{\mathbf{p} \in \Delta^{K-1}}{\operatorname{argmin}} \|\mathbf{p} - \mathbf{z}\|_2^2 \quad \text{No control over sparsity!!}$$



Sparsegen – Unified Framework

- A family of sparse probability mapping functions:

$$\rho(\mathbf{z}) = \text{sparsegen}(\mathbf{z}; g, \lambda) = \underset{\mathbf{p} \in \Delta^{K-1}}{\operatorname{argmin}} \|\mathbf{p} - g(\mathbf{z})\|_2^2 - \lambda \|\mathbf{p}\|_2^2 \quad g: \mathbb{R}^K \rightarrow \mathbb{R}^K$$

- Closed Form solution exists.

Key Contributions/Takeaways

- A unified framework for **sparse** probability mapping functions.
- Formulations *sparsegen-lin* and *sparsehourglass* – **control over sparsity**.
- Convex hinge-based loss functions** for multilabel classification.
- Sparser and more accurate prediction** for multilabel classification.
- Sparsity **control over attention heatmaps** in neural machine translation and abstractive summarization.

References

- [1] André F. T. Martins and Ramón F. Astudillo. (2016) *From softmax to sparsemax: A sparse model of attention and multi-label classification*. [ICML 2016]
- [2] Alexandre de Brébisson and Pascal Vincent. (2016) *An exploration of softmax alternatives belonging to the spherical loss family*. [ICLR 2016]
- [3] Vlad Niculae and Mathieu Blondel. (2017) *A regularized framework for sparse and structured neural attention*. [NIPS 2017]