# Linear models: Recap

Linear models:

- ▶ Perceptron

- ▶ Naïve Bay

$$\log P(y|\boldsymbol{x}; \boldsymbol{\theta}) = \log P(\boldsymbol{x}|y; \phi) + \mathsf{lo} \qquad B(\boldsymbol{x}) + \boldsymbol{\theta} \cdot \boldsymbol{f}(\boldsymbol{x}, y)$$

- ▶ Logistic Regression

$$\log P(y|\boldsymbol{x}; \boldsymbol{\theta}) = \boldsymbol{\theta} \cdot \boldsymbol{f}(\boldsymbol{x}, y) - \log \sum_{y' \in \mathcal{Y}} \exp \boldsymbol{\theta} \cdot \boldsymbol{f}(\boldsymbol{x}, y')$$

# Features and weights in linear models: Recap

- Feature represe

$$f(\boldsymbol{x}, y=1) = [\boldsymbol{x}; 0; 0; \underbrace{\cdots}; 0]$$
$$\phantom{f(\boldsymbol{x}, y=1) = [\boldsymbol{x}; 0; 0; \cdots; 0]}_{(K-1)}$$

$$f(\boldsymbol{x}, y=2) = [0; 0; \underbrace{\cdots}]$$
$$\phantom{xxxxxx}_{2) \times V}$$

$$\underbrace{\phantom{xxxxx}}_{(K}$$

- Weights: $\boldsymbol{\theta}$

$$\boldsymbol{\theta} = [\underbrace{\theta_1; \theta_2; \cdots; \theta_V}_{y=1}; \underbrace{\theta_1; \theta_2; \cdots; \theta_V}_{y=2}; \cdots; \underbrace{\theta_1; \theta_2; \cdots; \theta_V}_{y=K}]$$

# Rearranging the features and weights

▶ Represent the features $x$ as a *column* vector of length $V$, and represent the weights as a $\Theta$ as $K$

$$
x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_V \end{bmatrix}
\qquad
\Theta = \begin{array}{c} \\ y=2 \\ \vdots \\ y=K \end{array}
\begin{array}{cccc} x_1 & x_2 & \cdots & x_V \end{array}
\begin{bmatrix} & & & \theta_{1,V} \\ & & \cdot & \theta_{2,V} \\ & & \cdot & \\ \theta_{K,1} & \theta_{K,2} & \cdots & \theta_{K,V} \end{bmatrix}
$$

▶ What is $\Theta x$?

## Scores for each class

► Verify that $\psi_1, \psi_2, \ldots, \psi_K$ corres ach class

$$\Psi = \Theta x = $$

Implementation in Pytorch

https://eduassistpro.github.io/

Assignment Project Exam Help

Assignment Project Add WeChat edu_assist_pro

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# Digression: Matrix multiplication

▶ Matrix with $m$ rows and $n$ columns

$A \in \mathbb{R}^{m \times p}, B \in \mathbb{R}^{p \times n}$,

where $C_{ij}$

▶ Example:

$$\begin{bmatrix} 2 & 3 \\ 1 & 2 \end{bmatrix} \times \begin{bmatrix} 1 & 0 \\ 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} & & \\ 1 & 2 & 4 \end{bmatrix}$$

Digression: 3-D matrix multiplication

Tensor shape: (batch-size, sentence-length, embedding size)

# SoftMax

- SoftMax, also kn

$$\mathrm{SoftMax}_i(\psi) = \frac{\exp \psi_i}{\ \ \ \ \ }$$

for $i = 1, 2, \ldots, K$

- Applying Sof
distribution

$$\mathrm{SoftMax}(\boldsymbol{\Psi}) = \begin{bmatrix} P(y=1) \\ P(y=2) \\ \cdots \\ P(y=K) \end{bmatrix}$$

- Verify this is exactly logistic regression

Logistic regression as a neural network

$$\boldsymbol{y} = \mathsf{SoftMax}(\boldsymbol{\Theta}\boldsymbol{x})$$
$$V = 5 \, K = 3$$

# Going deep

► There is no reason w

$$z = \sigma(\boldsymbol{\Theta}_1 \boldsymbol{x})$$
$$\boldsymbol{y} = \text{SoftMax}(\boldsymbol{\Theta}_2 \boldsymbol{z})$$

# Going even deeper

- There is no reason w

$$z_1 = \sigma(\boldsymbol{\Theta}_1 \boldsymbol{x})$$
$$z_2 = \sigma(\boldsymbol{\Theta}_2 \boldsymbol{z}_1)$$
$$\boldsymbol{y} = \text{SoftMax}(\boldsymbol{\Theta}_3 \boldsymbol{z}_2)$$

- But why?

# Non-linear classification

Linear models like Logis[tic Regression represent documents in a]
high-dimensional vector space, and they are expressive enough and
work well for many NLP problems, why do we need more complex
non-linear models?

- ► There have been rapid advances in [deep learning], a family of
  nonlinear me[thods that transform the input] put
  through mul[tiple layers of nonlinear processing.]
- ► Deep learning facilitates the incorpora[tion of] **word**
  **embeddings**, which are dense vector r[epresentations of w]ords,
  that can be learned from massive amount[s of data.]
  - ► It has evolved from early static embeddings (e.g., Word2vec,
    Glove) to recent dynamtic embeddings (ELMO, BERT, XLNet)
- ► Rapid advances in specialized hardware called graphic
  processing units (GPUs). Many deep learning models can be
  implemented efficiently on GPUs.

# Feedforward Neural networks: an intuitive justification

- In image classification, instead of using the input (pixels) to predict the image type directly, you can imagi
  that you can predict the shapes or parts of an ima
  hand, ear.

- In text proce
  say we want to classify movie reviews (or m
  into a label set of {Good, Bad} ead predicting these
  labels directly, we first predict a set of comp
  such as the story, acting, soundtrack, cinematography, etc.
  from raw input (words in the text).

Face Recognition

https://eduassistpro.github.io/

Assignment Project Exam Help

Assignment Project Add WeChat edu_assist_pro

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

# Feedforward neural networks

Formally, this is what we do:

- **Use the text $x$ to predict the featur** ecifically, train a logistic regression classifier to compute or each $k \in \{1, 2, \cdots, K_z\}$

- **Use the featu** Train a logistic regression cl nknown or hidden, so we will use the $P(z|$ eatures.

Caveat: it's easy to demonstrate what this is wha for image processing, but it's hard to show this is what's actually going on in language processing. Interpretability is a major issue in neural models for language processing.

# The hidden layer: computing the composite features

- If we assume each $z_k$ is binary, that is, $z_k \in \{0, 1\}$, then $P(z_k|\boldsymbol{x})$ can be modeled with binary logistic regression

$$P(z_k = 1|\boldsymbol{x}; \boldsymbol{\Theta}^{(x \to z)}) = \sigma(\boldsymbol{\theta}^{x \to z}_k \cdot \boldsymbol{x}) = (1 + e^{-\boldsymbol{\theta}^{x \to z}_k \cdot \boldsymbol{x}})^{-1}$$

- The weight matrix is constructed by stacking (not concatenating, as in linear vectors for each $z_k$,

$$\boldsymbol{\Theta}^{(x \to z)} = \left[\boldsymbol{\theta}^{x \to z}_1, \boldsymbol{\theta}^{x \to z}_2, \cdots, \boldsymbol{\theta}^{x \to}_{K_z}\right.$$

- We assume an offset/bias term is included in $\boldsymbol{x}$ and its parameter is included in each $\boldsymbol{\theta}^{x \to z}_k$

Notations: $\boldsymbol{\Theta}^{(x \to z)} \in \mathbb{R}^{k_z \times V}$ is a real number matrix with a dimension of $k_z$ rows and $V$ columns

# The output layer

▶ The output layer is computed by the multiclass logistic regression probability

$$P(y = j \mid \boldsymbol{z}; \boldsymbol{\Theta}^{(z \to y)}, \boldsymbol{b}) = \frac{\exp(\boldsymbol{\theta}_j^{(z \to y)} \cdot \boldsymbol{z} + b_j)}{\sum \exp(\boldsymbol{\theta}^{(z \to y)} \cdot \boldsymbol{z} + b_j')}$$

▶ The weight matrix $\boldsymbol{\Theta}^{(z \to y)} \in \mathbb{R}^{K_y \times K_z}$ is constructed by stacking weight vectors for each

$$\boldsymbol{\Theta}^{(z \to y)} = \left[ \boldsymbol{\theta}_1^{z \to y}, \boldsymbol{\theta}_2^{z \to y}, \cdots, \boldsymbol{\theta}_{K_y}^{z} \right]$$

▶ The vector of probabilities over each possible value of $y$ is denoted:
$$P(\boldsymbol{y} \mid \boldsymbol{z}; \boldsymbol{\Theta}^{(z \to y)}, \boldsymbol{b}) = \mathrm{SoftMax}(\boldsymbol{\Theta}^{(z \to y)} \boldsymbol{z} + \boldsymbol{b})$$

# Activation functions

▶ Sigmoid: The

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

▶ Tanh: The range of the tanh activation function is $(-1, 1)$

▶ ReLU: The **rectified linear unit (ReLU)** or negative inputs, and linear for positive inputs

$$ReLU(x) = max(x, 0) = \begin{cases} 0 & x < 0 \\ x & otherwise \end{cases}$$

Sigmoid and tanh are sometimes described as **squashing functions**.

# Activation functions in Pytorch

```python
from torch import nn
import torch

input = torch.ra     4)
sigmoid = nn.Si
output = sigmoid

tanh = nn.Tanh()
output = tanh(input)

relu = nn.ReLU()
output = relu(input)
```

# Output and loss functions

In a multi-class classific
probabilistic distribution over possible labels. It works well together
with negative conditional likelihood (just like logistic regression)

$$-\mathcal{L} = -\sum_{i=1}^{N} \log P(y^{(i}$$

or **cross entropy l**

$$\tilde{y}_j \triangleq Pr(y =$$

$$-\mathcal{L} = -\sum_{i=1}^{N} \boldsymbol{e}_{y^{(i)}} \cdot \log \tilde{\boldsymbol{y}}$$

where $\boldsymbol{e}_{y^{(i)}}$ is a **one-hot vector** of zeros with a value of one at the position $y^{(i)}$

# Output and loss function

- There are alternatives to SoftMax and cross-entropy loss, just as there are alternatives in linear models.
- Pairing an affine transformation (rememb               h a margin loss:

$$\Psi(y; \boldsymbol{x}^{(i)}, \Theta) = \boldsymbol{\theta}_y^{(z \to y)} \cdot \boldsymbol{z} + b_y$$

$$\ell_{\mathsf{MARGIN}}(\Theta; \boldsymbol{x}^{(i)}, y^{(i)}) = \max_{y \neq y^{(i)}} \left(1 + \Psi(y; \boldsymbol{x}^{(i)}, \Theta) - \Psi(y^{(i)}; \boldsymbol{x}^{(i)}, \Theta)\right)_+$$

# Inputs and Lookup layers

- Assuming a bag-of-words model, when the input $x$ is the count of each word $x_j$. (This can be generalized to feature count).

- To compute the hidden unit $z_k$:

- This text representation is particularly s̶ ard networks.

- The connections from word $j$ to each of the hidden units $z_k$ form a vector $\theta_j^{(x \to z)}$ is sometimes described as the embedding of word $j$. Word embeddings can be learned from unlabeled data, using techniques such as Word2Vec and GLOVE.

# Alternative text representations

- ▶ Alternatively, a text can be represented as a sequence of word tokens $w_1, w_2, w_3, \cdots, w_M$. This view is useful for models such as **Convolutional Neural Netw** **vNets**, which processes text as a sequence.

- ▶ Each word tok tor $\boldsymbol{e}_{w_m}$, with dimensi represented e-hot vectors: $\boldsymbol{W} = [\boldsymbol{e}_{w_1}, \boldsymbol{e}_{w_2}, \cdots, \boldsymbol{e}_{w}$

- ▶ To show that this is equivalent to the bag-o can recover the word count from the matrix-vector product $\boldsymbol{W}[1, 1, \cdots, 1]^\top \in R^V$.

- ▶ The matrix product $\Theta^{x \to z} \boldsymbol{W} \in R^{k_z \times M}$ contains horizontally concatenated embeddings of each word in the document.