

EE 628

Deep Learning

Fall 2019

Lecture 2
09/05/2019

Assistant Prof. Sergul Aydore
Department of Electrical and Computer Engineering



Linear Neural Networks

- We need to cover the basics of neural network training

Linear Neural Networks

- We need to cover the basics of neural network training
- Entire training process includes

Linear Neural Networks

- We need to cover the basics of neural network training
- Entire training process includes
 - Defining simple neural network architectures
 - Handling data
 - Specifying loss functions
 - Training the model

Linear Regression

- **Linear:** the output is a linear combination of the input features

Linear Regression

- **Linear:** the output is a linear combination of the input features
- Example: estimating a price of a house (in dollars) based on area (in square feet) and age (in years)

Linear Regression

- **Linear:** the output is a linear combination of the input features
- Example: estimating a price of a house (in dollars) based on area (in square feet) and age (in years)
 - Our model can be expressed as:
 $\text{price} = w_{\text{area}} \times \text{area} + w_{\text{age}} \times \text{age} + b$

Linear Regression

- **Linear:** the output is a linear combination of the input features
- Example: estimating a price of a house (in dollars) based on area (in square feet) and age (in years)
 - Our model can be expressed as:
$$\text{price} = w_{\text{area}} \times \text{area} + w_{\text{age}} \times \text{age} + b$$
- In the case of d variables: $\hat{y} = w_1 \times x_1 + \cdots + w_d \times x_d + b$
 - In vector form:

$$\hat{y} = \mathbf{w}^\top \mathbf{x} + b$$

Linear Regression

- **Linear:** the output is a linear combination of the input features
- Example: estimating a price of a house (in dollars) based on area (in square feet) and age (in years)
 - Our model can be expressed as:
$$\text{price} = w_{\text{area}} \times \text{area} + w_{\text{age}} \times \text{age} + b$$
- In the case of d variables: $\hat{y} = w_1 \times x_1 + \cdots + w_d \times x_d + b$
 - In vector form:
$$\hat{y} = \mathbf{w}^T \mathbf{x} + b$$
- For a collection of data points \mathbf{X} , the prediction can be expressed as:

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{w} + \mathbf{b}$$



Design matrix

Training Data

- We need a collection of examples for which we know both the actual selling price of each house as well as their corresponding area and age.

Training Data

- We need a collection of examples for which we know both the actual selling price of each house as well as their corresponding area and age.
- Goal: minimize the error between the predicted price and the real price

Training Data

- We need a collection of examples for which we know both the actual selling price of each house as well as their corresponding area and age.
- Goal: minimize the error between the predicted price and the real price
- A house comprises one *sample*:
 - Its actual selling price is called a *label*
 - The factors used to predict the price are called *features*

Training Data

- We need a collection of examples for which we know both the actual selling price of each house as well as their corresponding area and age.
- Goal: minimize the error between the predicted price and the real price
- A house comprises one *sample*:
 - Its actual selling price is called a *label*
 - The factors used to predict the price are called *features*
- n : number of samples in our dataset

Training Data

- We need a collection of examples for which we know both the actual selling price of each house as well as their corresponding area and age.
- Goal: minimize the error between the predicted price and the real price
- A house comprises one *sample*:
 - Its actual selling price is called a *label*
 - The factors used to predict the price are called *features*
- n : number of samples in our dataset
- We index the samples by i :
 - Each input data point: $x^{(i)} = [x_1^{(i)}, x_2^{(i)}]$
 - The corresponding label: $y^{(i)}$

Loss Function

- Need to measure error between the predicted value and the real value of the price: i.e. *loss function*

Loss Function

- Need to measure error between the predicted value and the real value of the price: i.e. *loss function*
- The smaller the value, the smaller the error

Loss Function

- Need to measure error between the predicted value and the real value of the price: i.e. *loss function*
- The smaller the value, the smaller the error
- A common choice: the *square loss*

Loss Function

- Need to measure error between the predicted value and the real value of the price: i.e. *loss function*
- The smaller the value, the smaller the error
- A common choice: the *square loss*
- For given parameters \mathbf{w} and b , we can express the error:

$$l^{(i)}(\mathbf{w}, b) = \frac{1}{2} (\hat{y}^{(i)} - y^{(i)})^2$$

Loss Function

- Need to measure error between the predicted value and the real value of the price: i.e. *loss function*
- The smaller the value, the smaller the error
- A common choice: the *square loss*
- For given parameters \mathbf{w} and b , we can express the error:

$$l^{(i)}(\mathbf{w}, b) = \frac{1}{2} (\hat{y}^{(i)} - y^{(i)})^2$$

- To measure the quality of a model on the entire dataset, we average the losses on the training set.

$$L(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n l^{(i)}(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (\mathbf{w}^\top \mathbf{x}^{(i)} + b - y^{(i)})^2$$

Loss Function

- Need to measure error between the predicted value and the real value of the price: i.e. *loss function*
- The smaller the value, the smaller the error
- A common choice: the *square loss*
- For given parameters \mathbf{w} and b , we can express the error:

$$l^{(i)}(\mathbf{w}, b) = \frac{1}{2} (\hat{y}^{(i)} - y^{(i)})^2$$

- To measure the quality of a model on the entire dataset, we average the losses on the training set.

$$L(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n l^{(i)}(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (\mathbf{w}^\top \mathbf{x}^{(i)} + b - y^{(i)})^2$$

- We want to find parameters that minimize the average loss across all training samples

$$\mathbf{w}^*, b^* = \arg \min_{\mathbf{w}, b} L(\mathbf{w}, b)$$

Analytic Solution

- Linear regression is unusually a simple optimization problem

Analytic Solution

- Linear regression is unusually a simple optimization problem
- There is just one critical point on the loss surface corresponding to the global minimum.

Analytic Solution

- Linear regression is unusually a simple optimization problem
- There is just one critical point on the loss surface corresponding to the global minimum.
- Taking the derivative of the loss with respect to \mathbf{w} , and setting it equal to 0 gives the analytic solution:

$$\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

Gradient Descent

- The key trick behind nearly all deep learning is gradient descent:
 - To reduce the error gradually by iteratively updating the parameters

Follow the slope!



Credit: http://cs231n.stanford.edu/slides/2018/cs231n_2018_lecture03.pdf

Gradient Descent

- The key trick behind nearly all deep learning is gradient descent:
 - To reduce the error gradually by iteratively updating the parameters
 - On convex loss surfaces it will eventually converge to a global minimum

Follow the slope!



Credit: http://cs231n.stanford.edu/slides/2018/cs231n_2018_lecture03.pdf

Gradient Descent

- The key trick behind nearly all deep learning is gradient descent:
 - To reduce the error gradually by iteratively updating the parameters
- On convex loss surfaces it will eventually converge to a global minimum
- For nonconvex surfaces, it will at least lead towards a (hopefully good) local minimum

Follow the slope!



Credit: http://cs231n.stanford.edu/slides/2018/cs231n_2018_lecture03.pdf

The Normal Distribution and Squared Loss

- Recall the normal distribution with mean μ and variance σ^2 : $p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$

The Normal Distribution and Squared Loss

- Recall the normal distribution with mean μ and variance σ^2 : $p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$
- The key assumption in linear regression with least mean squares loss is that the observations actually arise from noisy observations: $y = \mathbf{w}^\top \mathbf{x} + b + \epsilon$ where $\epsilon \sim \mathcal{N}(0, \sigma^2)$

The Normal Distribution and Squared Loss

- Recall the normal distribution with mean μ and variance σ^2 : $p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$
- The key assumption in linear regression with least mean squares loss is that the observations actually arise from noisy observations: $y = \mathbf{w}^\top \mathbf{x} + b + \epsilon$ where $\epsilon \sim \mathcal{N}(0, \sigma^2)$
- This allows us to write out the likelihood of seeing a particular y for a given data \mathbf{x} :

$$p(y|\mathbf{x}) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{1}{2\sigma^2}(y - \mathbf{w}^\top \mathbf{x} - b)^2\right)$$

The Normal Distribution and Squared Loss

- Recall the normal distribution with mean μ and variance σ^2 : $p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$
- The key assumption in linear regression with least mean squares loss is that the observations actually arise from noisy observations: $y = \mathbf{w}^\top \mathbf{x} + b + \epsilon$ where $\epsilon \sim \mathcal{N}(0, \sigma^2)$

- This allows us to write out the likelihood of seeing a particular y for a given data \mathbf{x} :

$$p(y|\mathbf{x}) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{1}{2\sigma^2}(y - \mathbf{w}^\top \mathbf{x} - b)^2\right)$$

- A good way of finding the most likely values of b and w is to maximize the likelihood of the entire dataset $p(\mathbf{y}|\mathbf{x}) = \prod_{i=1} p(y^{(i)}|\mathbf{x}^{(i)})$

The Normal Distribution and Squared Loss

- Recall the normal distribution with mean μ and variance σ^2 : $p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$
- The key assumption in linear regression with least mean squares loss is that the observations actually arise from noisy observations: $y = \mathbf{w}^\top \mathbf{x} + b + \epsilon$ where $\epsilon \sim \mathcal{N}(0, \sigma^2)$
- This allows us to write out the likelihood of seeing a particular y for a given data \mathbf{x} :

$$p(y|\mathbf{x}) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{1}{2\sigma^2}(y - \mathbf{w}^\top \mathbf{x} - b)^2\right)$$

- A good way of finding the most likely values of b and w is to maximize the likelihood of the entire dataset $p(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^n p(y^{(i)}|\mathbf{x}^{(i)})$
- The estimators found by maximizing the likelihood of the data is known as *Maximum likelihood estimators*.

The Normal Distribution and Squared Loss

- Recall the normal distribution with mean μ and variance σ^2 : $p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$
- The key assumption in linear regression with least mean squares loss is that the observations actually arise from noisy observations: $y = \mathbf{w}^\top \mathbf{x} + b + \epsilon$ where $\epsilon \sim \mathcal{N}(0, \sigma^2)$
- This allows us to write out the likelihood of seeing a particular y for a given data \mathbf{x} :

$$p(y|\mathbf{x}) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{1}{2\sigma^2}(y - \mathbf{w}^\top \mathbf{x} - b)^2\right)$$

- A good way of finding the most likely values of b and w is to maximize the likelihood of the entire dataset $p(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^n p(y^{(i)}|\mathbf{x}^{(i)})$
- The estimators found by maximizing the likelihood of the data is known as *Maximum likelihood estimators*.
- Maximizing the product of many exponential functions is awkward. A much better way is to minimize the Negative Log-likelihood- $\log p(\mathbf{y} | \mathbf{x})$:

$$-\log(\mathbf{y}|\mathbf{x}) = \sum_{i=1}^n \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2}(y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)} - b)^2$$

The Normal Distribution and Squared Loss

- Recall the normal distribution with mean μ and variance σ^2 : $p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$
- The key assumption in linear regression with least mean squares loss is that the observations actually arise from noisy observations: $y = \mathbf{w}^\top \mathbf{x} + b + \epsilon$ where $\epsilon \sim \mathcal{N}(0, \sigma^2)$
- This allows us to write out the likelihood of seeing a particular y for a given data \mathbf{x} :

$$p(y|\mathbf{x}) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{1}{2\sigma^2}(y - \mathbf{w}^\top \mathbf{x} - b)^2\right)$$

- A good way of finding the most likely values of b and w is to maximize the likelihood of the entire dataset $p(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^n p(y^{(i)}|\mathbf{x}^{(i)})$
- The estimators found by maximizing the likelihood of the data is known as *Maximum likelihood estimators*.
- Maximizing the product of many exponential functions is awkward. A much better way is to minimize the Negative Log-likelihood- $\log p(\mathbf{y} | \mathbf{x})$:

$$-\log(\mathbf{y}|\mathbf{x}) = \sum_{i=1}^n \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2}(y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)} - b)^2$$

- It follows that maximum likelihood in a linear model with additive Gaussian noise is equivalent to linear regression with squared loss

Implementation of linear regression from scratch

- Now, start your ipython notebook

Softmax Regression

- In practice, sometimes we are interested in classification: asking not *how much* but *which one*.

Softmax Regression

- In practice, sometimes we are interested in classification: asking not *how much* but *which one*.
 - Does this email belong in the spam folder or the inbox?

Softmax Regression

- In practice, sometimes we are interested in classification: asking not *how much* but *which one*.
 - Does this email belong in the spam folder or the inbox?
 - Is this customer more likely to sign up or not to sign up for a subscription service?

Softmax Regression

- In practice, sometimes we are interested in classification: asking not *how much* but *which one*.
 - Does this email belong in the spam folder or the inbox?
 - Is this customer more likely to sign up or not to sign up for a subscription service?
 - Does this image depict a donkey, a dog, a cat, or a rooster?

Softmax Regression

- In practice, sometimes we are interested in classification: asking not *how much* but *which one*.
 - Does this email belong in the spam folder or the inbox?
 - Is this customer more likely to sign up or not to sign up for a subscription service?
 - Does this image depict a donkey, a dog, a cat, or a rooster?
 - Which movie is user most likely to watch next?

Classification Problems

- Let's start with a really simple image classification problem

Classification Problems

- Let's start with a really simple image classification problem
- Each input will be a grayscale 2×2 image

Classification Problems

- Let's start with a really simple image classification problem
- Each input will be a grayscale 2×2 image
- Each image can be represented with 4 features: x_1, x_2, x_3, x_4

Classification Problems

- Let's start with a really simple image classification problem
- Each input will be a grayscale 2×2 image
- Each image can be represented with 4 features: x_1, x_2, x_3, x_4
- Let's assume each image belongs to one among three categories: *cat*, *chicken* and *dog*

Classification Problems

- Let's start with a really simple image classification problem
- Each input will be a grayscale 2×2 image
- Each image can be represented with 4 features: x_1, x_2, x_3, x_4
- Let's assume each image belongs to one among three categories: *cat*, *chicken* and *dog*
- We need to represent the labels. Which one makes more sense?
 - $y \in \{1, 2, 3\}$
 - $y \in \{(1, 0, 0), (0, 1, 0), (0, 0, 1)\}$

Network Architecture

- In order to estimate multiple classes, we need a model with multiple outputs, one per category.

Network Architecture

- In order to estimate multiple classes, we need a model with multiple outputs, one per category.
- We compute 3 outputs for each input

$$o_1 = x_1 w_{11} + x_2 w_{21} + x_3 w_{31} + x_4 w_{41} + b_1$$

$$o_2 = x_1 w_{12} + x_2 w_{22} + x_3 w_{32} + x_4 w_{42} + b_2$$

$$o_3 = x_1 w_{13} + x_2 w_{23} + x_3 w_{33} + x_4 w_{43} + b_3$$

— — —

Network Architecture

- In order to estimate multiple classes, we need a model with multiple outputs, one per category.
- We compute 3 outputs for each input

$$o_1 = x_1 w_{11} + x_2 w_{21} + x_3 w_{31} + x_4 w_{41} + b_1$$

$$o_2 = x_1 w_{12} + x_2 w_{22} + x_3 w_{32} + x_4 w_{42} + b_2$$

$$o_3 = x_1 w_{13} + x_2 w_{23} + x_3 w_{33} + x_4 w_{43} + b_3$$

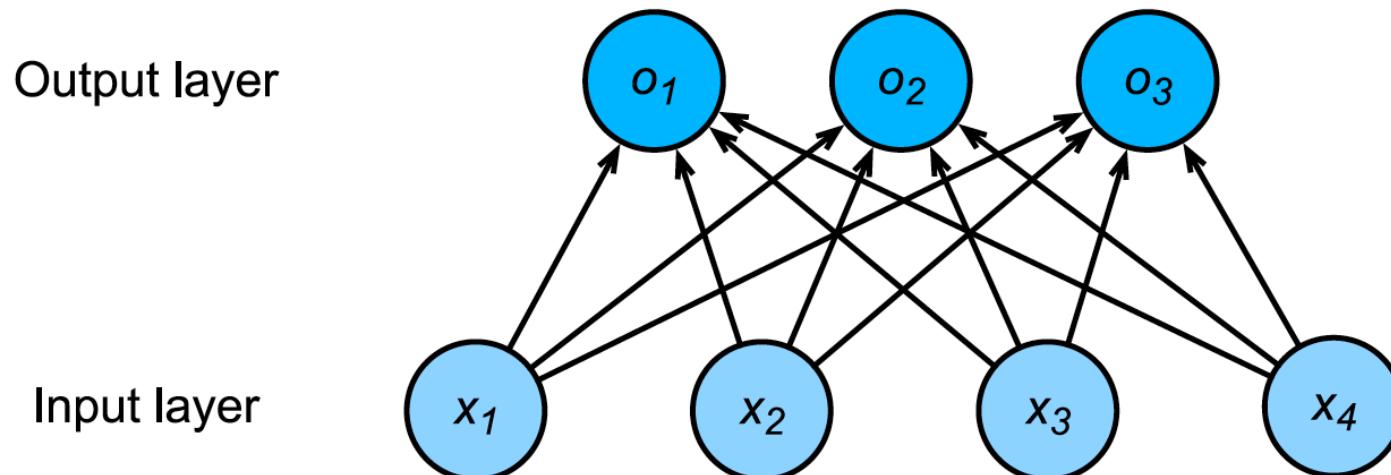


Fig. 5.4.1: Softmax regression is a single-layer neural network.

Softmax Operation

- Let's use linear algebra notation for our model $\mathbf{o} = \mathbf{W}\mathbf{x} + \mathbf{b}$

Softmax Operation

- Let's use linear algebra notation for our model $\mathbf{o} = \mathbf{W}\mathbf{x} + \mathbf{b}$
- We want our outputs to correspond to probabilities to indicate our confidence that the item belongs to each category

Softmax Operation

- Let's use linear algebra notation for our model $\mathbf{o} = \mathbf{W}\mathbf{x} + \mathbf{b}$
- We want our outputs to correspond to probabilities to indicate our confidence that the item belongs to each category
- Challenges:
 - We need to guarantee that each output is nonnegative
 - Probabilities sum up to 1

Softmax Operation

- Let's use linear algebra notation for our model $\mathbf{o} = \mathbf{W}\mathbf{x} + \mathbf{b}$
- We want our outputs to correspond to probabilities to indicate our confidence that the item belongs to each category
- Challenges:
 - We need to guarantee that each output is nonnegative
 - Probabilities sum up to 1
- Softmax function does it all!

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{o}) \text{ where } \hat{y}_i = \frac{\exp(o_i)}{\sum_j \exp(o_j)}$$

Softmax Operation

- Let's use linear algebra notation for our model $\mathbf{o} = \mathbf{W}\mathbf{x} + \mathbf{b}$
- We want our outputs to correspond to probabilities to indicate our confidence that the item belongs to each category
- Challenges:
 - We need to guarantee that each output is nonnegative
 - Probabilities sum up to 1
- Softmax function does it all!

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{o}) \text{ where } \hat{y}_i = \frac{\exp(o_i)}{\sum_j \exp(o_j)}$$

- Vectorization for minibatches:

$$\mathbf{O} = \mathbf{X}\mathbf{W} + \mathbf{b}$$

$$\hat{\mathbf{Y}} = \text{softmax}(\mathbf{O})$$

Loss Function

- Log likelihood: check how well we predicted what we observe

$$p(Y|X) = \prod_{i=1}^n p(y^{(i)}|x^{(i)}) \text{ and thus } -\log p(Y|X) = \sum_{i=1}^n -\log p(y^{(i)}|x^{(i)})$$
$$-\log p(y^{(i)}|x^{(i)}) = -\sum y_j \log \hat{y}_j$$

Loss Function

- Log likelihood: check how well we predicted what we observe

$$p(Y|X) = \prod_{i=1}^n p(y^{(i)}|x^{(i)}) \text{ and thus } -\log p(Y|X) = \sum_{i=1}^n -\log p(y^{(i)}|x^{(i)})$$

$$-\log p(y^{(i)}|x^{(i)}) = -\sum_j y_j \log \hat{y}_j$$

- Compute the derivative with respect to o_j .

Loss Function

- Log likelihood: check how well we predicted what we observe

$$p(Y|X) = \prod_{i=1}^n p(y^{(i)}|x^{(i)}) \text{ and thus } -\log p(Y|X) = \sum_{i=1}^n -\log p(y^{(i)}|x^{(i)})$$
$$-\log p(y^{(i)}|x^{(i)}) = -\sum_j y_j \log \hat{y}_j$$

- Compute the derivative with respect to o_j .
- Cross entropy loss: one of the most commonly used losses for multiclass classification

$$l(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_j y_j \log(\hat{y}_j)$$

Loss Function

- Log likelihood: check how well we predicted what we observe

$$p(Y|X) = \prod_{i=1}^n p(y^{(i)}|x^{(i)}) \text{ and thus } -\log p(Y|X) = \sum_{i=1}^n -\log p(y^{(i)}|x^{(i)})$$
$$-\log p(y^{(i)}|x^{(i)}) = -\sum_j y_j \log \hat{y}_j$$

- Compute the derivative with respect to o_j .
- Cross entropy loss: one of the most commonly used losses for multiclass classification $l(\mathbf{y}, \hat{\mathbf{y}}) = -\sum y_j \log(\hat{y}_j)$
- Kullback Leibler Divergence: measures similarity between two distributions

$$D(p||q) = \sum_j p(j) \log \frac{p(j)}{q(j)}$$

Loss Function

- Log likelihood: check how well we predicted what we observe

$$p(Y|X) = \prod_{i=1}^n p(y^{(i)}|x^{(i)}) \text{ and thus } -\log p(Y|X) = \sum_{i=1}^n -\log p(y^{(i)}|x^{(i)})$$
$$-\log p(y^{(i)}|x^{(i)}) = -\sum_j y_j \log \hat{y}_j$$

- Compute the derivative with respect to o_j .
- Cross entropy loss: one of the most commonly used losses for multiclass classification $l(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_j y_j \log(\hat{y}_j)$
- Kullback Leibler Divergence: measures similarity between two distributions

$$D(p||q) = \sum_j p(j) \log \frac{p(j)}{q(j)}$$

- Minimizing $D(p || q)$ with respect to q is equivalent to minimizing the cross-entropy loss. (PROVE!)

Next

- Multilayer perceptron