

Machine Learning.

Introduction to convolutional neural networks.

Example. LeNet Mathematica.

Definition. ~~It is~~ We define an order K tensor to be an ~~extended vector~~ K -dimensional array of numbers $X \in \mathbb{R}^{d_1 \times \dots \times d_K}$

Examples.

Vectors are order 1 tensors $X \in \mathbb{R}^n$

Matrices are order 2 tensors $X \in \mathbb{R}^{m \times n}$

An RGB image which is H pixels high and W pixels wide is represented by an ~~order~~ order 3 tensor

$$X \in \mathbb{R}^{H \times W \times 4}$$

The elements of an order K tensor are denoted with K indices: $X_{i_1 \dots i_K} \in \mathbb{R}$

Definition. The operator vec reduces every tensor to a vector by ordering the entries in reverse lexicographic order.

Example

$$\text{vec} \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \\ a_{31} & a_{32} \end{bmatrix} = \begin{bmatrix} a_{11} \\ a_{21} \\ a_{31} \\ a_{12} \\ a_{22} \\ a_{32} \end{bmatrix}$$

$$\text{vec} \begin{array}{c} \begin{array}{ccc} a_{111} & a_{112} & a_{122} \\ & a_{21} & \\ a_{121} & & a_{222} \end{array} \\ \begin{array}{cc} a_{211} & a_{221} \end{array} \end{array} = \begin{bmatrix} a_{111} \\ a_{211} \\ a_{121} \\ a_{221} \\ a_{112} \\ a_{212} \\ a_{122} \\ a_{222} \end{bmatrix}$$

Recall.

Suppose $z: \mathbb{R}^m \rightarrow \mathbb{R}$, and we write $z = z(\vec{y})$.

(3)

Then

$$\frac{\partial z}{\partial \vec{y}^T} = Dz, \text{ the } \begin{matrix} \text{1-tensor} \\ \text{matrix} \end{matrix} \text{ so that } \left[\frac{\partial z}{\partial \vec{y}} \right]_i = \frac{\partial z}{\partial y_i}$$

is a $1 \times m$ matrix of partial derivatives.

If $\vec{y}: \mathbb{R}^n \rightarrow \mathbb{R}^m$, and we write $\vec{y} = \vec{y}(\vec{x})$
and

$$\frac{\partial \vec{y}}{\partial \vec{x}^T} = D\vec{y}, \text{ the } \begin{matrix} \text{2-tensor} \\ \text{matrix} \end{matrix} \text{ so } \left[\frac{\partial \vec{y}}{\partial \vec{x}^T} \right]_{ij} = \frac{\partial y_i}{\partial x_j}$$

Then we can write the chain rule as

$$\frac{\partial z}{\partial \vec{x}^T} = \frac{\partial z}{\partial \vec{y}^T} \frac{\partial \vec{y}}{\partial \vec{x}^T} \text{ or } Dz Dy$$

↑ ↑ ↖ $m \times n$ matrix

$1 \times n$ row vector row vector $1 \times m$

Definition. A CNN is a composition of maps (called layers) defined by tensors $\overline{w}_1, \dots, \overline{w}_L$ which we denote $\boxed{\overline{w}_i}$. Each map takes tensors to tensors so that the composition

$$X^1 \rightarrow \boxed{\overline{w}^1} \rightarrow X^2 \rightarrow \boxed{\overline{w}^2} \rightarrow \dots \rightarrow X^L \rightarrow \boxed{\overline{w}^L} \rightarrow \mathbb{Z}$$

is defined.

To classify images ~~we~~ into C categories,

X^1 = input image order 3, $H \times W \times 3$ and X^L = output prediction
order 1, C

X^1_{ijk} = color k of
pixel ij

X^L_i = probability
that image
is category i

(5)

The last layer is special. If we know the classification of image X^L is category i , then (for instance)

$$z = \frac{1}{2} \|\vec{e}_i - X^L\|^2$$

The choice of function here is part of the model (called a loss function).

Given the parameters W^1, \dots, W^{L-1} , the CNN prediction ~~is given by~~ for image X^1 is given by ~~$\arg \max$~~ $\arg \max_{i \in 1, \dots, C} X_i^L = i \mid X_i^L \geq X_j^L$ for all $j \in 1, \dots, C$

Idea. Given a "training set" of N images, the total loss is a function z of W^1, \dots, W^{L-1} . ~~z~~ We want to minimize

$$z(W^1, \dots, W^{L-1})$$

6

Note: To define the loss layer, we need to know the truth for each image in the training set.

Note 2: The w^1, \dots, w^{L-1} are a lot of variables, but in theory we really just want to solve

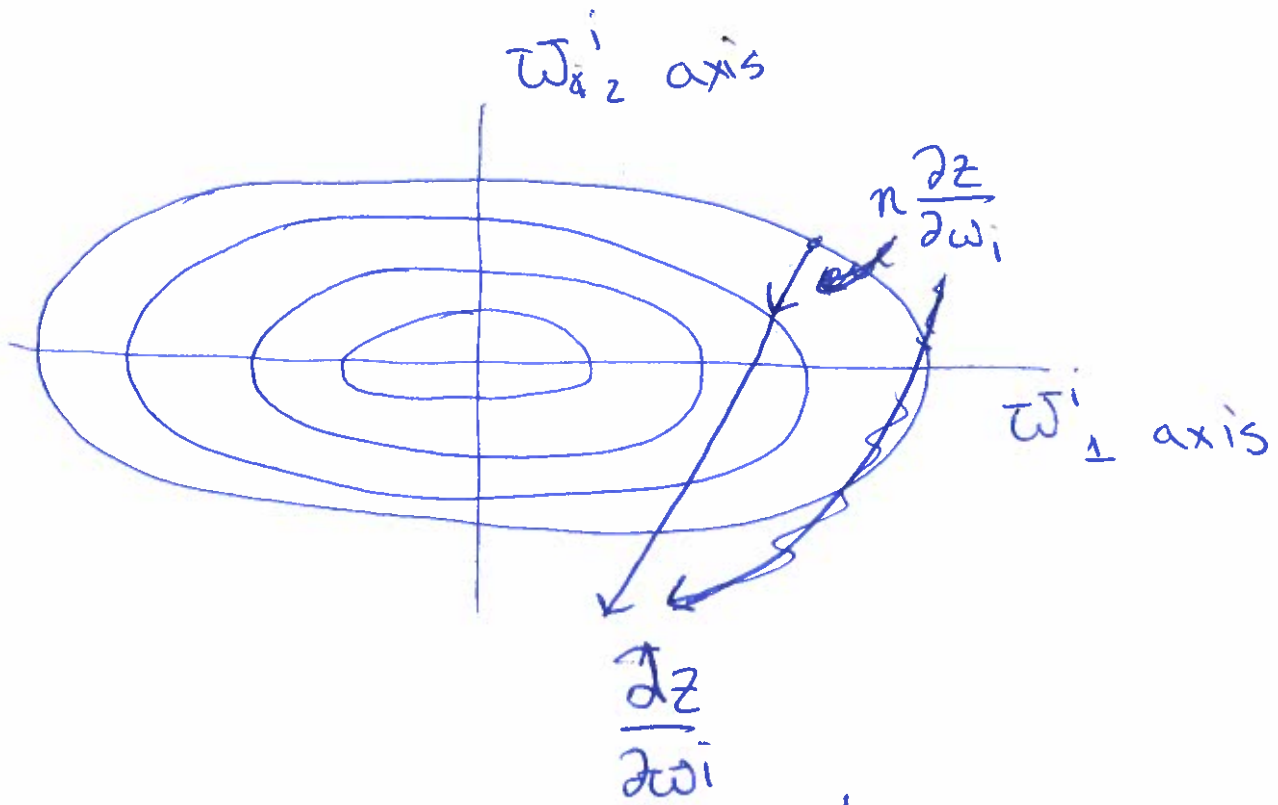
$$\frac{\partial z}{\partial (w^1, \dots, w^{L-1})} = \vec{0}$$

We can update w^i by steepest descent by the rule

$$w^i \leftarrow w^i - \eta \underbrace{\frac{\partial z}{\partial w^i}}$$

a small number called the learning rate.
 a tensor of the same dimensions as w^i

7



Problem. If we minimize ^{loss} on a single example, we will get that example right but others wrong. If we minimize on all examples, # computations is large.

Solution. Choose ~~2~~ ^K ~~or 4~~ ^{randomly} examples, let X^1 be $H \times W \times 3 \times$ ~~2~~ K and compute a gradient step for w^1, \dots, w^{L-1} , then ~~choose~~ & repeat.

But how do we compute the $\frac{\partial z}{\partial w_i}$?

Back propagation. For every layer, compute

$$\frac{\partial z}{\partial x^i} \text{ and } \frac{\partial z}{\partial w_i}$$

Start with the last layer $\boxed{w^L}$ which has no parameters, so

$$\frac{\partial z}{\partial w^L} = 0, \quad \frac{\partial z}{\partial (x^L)^k} = \overline{x^L} - t$$

↑
truth.