THzürich



Background Review

Reliable and Interpretable Artificial Intelligence

Assumed Background

Logic

- $\bullet \ \land, \lor, \rightarrow, \Rightarrow$
- ∀,∃
- Predicates

Linear Algebra

- Vectors
- Matrices

Probability Theory

- Random variables, Indicator variables
- Probabilities
- Bayes' law
- Expectation
- Distributions
 - Cumulative distribution function (CDF)
 - Probability density function (PDF)

Linear Algebra

Linear transformations

Definition

A linear transformation is a map $f : \mathbb{R}^m \to \mathbb{R}^n$ that preserves linear combinations:

$$f\left(\sum_{i=1}^{k}\lambda_{i}\vec{u}_{i}\right)=\sum_{i=1}^{k}\lambda_{i}f\left(\vec{u}_{i}\right).$$

Example (in
$$\mathbb{R}^2$$
)

Linear transformations as matrices

Proposition

Every linear transformation $f : \mathbb{R}^m \to \mathbb{R}^n$ can be expressed as a matrix $M \in \mathbb{R}^{n \times m}$:

 $\forall \vec{u} \in \mathbb{R}^m \colon f(\vec{u}) = M\vec{u}.$

Example (in \mathbb{R}^2)

$$\begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$$

Linear transformations as matrices

Proposition

Every linear transformation $f : \mathbb{R}^m \to \mathbb{R}^n$ can be expressed as a matrix $M \in \mathbb{R}^{n \times m}$:

 $\forall \vec{u} \in \mathbb{R}^m \colon f(\vec{u}) = M\vec{u}.$

Example (in \mathbb{R}^2)

$$\begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$$
rotation dilation projection

Affine transformations

Definition

An affine transformation is a map $f : \mathbb{R}^m \to \mathbb{R}^n$ that preserves affine combinations:

$$\sum_{i=1}^{k} \lambda_{i} = 1 \implies f\left(\sum_{i=1}^{k} \lambda_{i} \vec{u}_{i}\right) = \sum_{i=1}^{k} \lambda_{i} f\left(\vec{u}_{i}\right).$$

Proposition

Every affine $f : \mathbb{R}^m \to \mathbb{R}^n$ *decomposes into a translation after a linear transformation:*

$$f(\vec{u})=\vec{t}+g(\vec{u}).$$

ETH zürich

Norms and distances

Euclidean (*I*₂)

$$\|\vec{u}\|_2 = \sqrt[2]{\vec{u} \cdot \vec{u}}$$

$$d_2(\vec{u}, \vec{v}) = \|\vec{v} - \vec{u}\|_2.$$

 $l_{
ho}$ (1 \leq ho \leq ∞)

$$\|\vec{u}\|_p = \sqrt[p]{\sum |u_i|^p}$$

$$d_p(\vec{u},\vec{v}) = \|\vec{v} - \vec{u}\|_p.$$

$$\|\vec{u}\|_{\infty} = \lim_{p \to \infty} \|\vec{u}\|_p = \sup |u_i|.$$

DINFK

 I_{∞}

Prediction

Prediction problems

Predict lifespan given GDP per capita.



Predict quantity \implies regression.

Predict digit given a pixelated scan.



MNIST http://yann.lecun.com/exdb/mnist/

Predict label \implies classification.

Mathematical formulation

The simple version

Predict a latent variables \vec{Y} from a manifest variables \vec{X} , i.e., select a model f such that $f(\vec{X}) \approx \vec{Y}$.

Consistency

- 1. In practice, we need to sample (\vec{X}, \vec{Y}) in order to select *f*.
- 2. This means that $f = f_n$ depends on the sample size *n*.
- 3. The selection method is consistent if $f_n(\vec{X}) \rightarrow \vec{Y}$ as $n \rightarrow \infty$.

The not so simple version

What is $f(\vec{X}) \approx \vec{Y}$?

- 1. No universal answer.
- 2. A choice depending on the task.
- 3. Usually defined by a loss function:

 $\mathcal{L}_f(\vec{x}, \vec{y}) \in \mathbb{R}$

Where does *f* come from?

- 1. Comes from a model space F.
- 2. That space is, again, a choice.
- 3. Usually parameterized by a vector:

$$F = \{f_{\theta} \mid \theta \in \mathbb{R}^d\}.$$

Model fitting

Ideally, we want a model $f_{\theta^*} \in F$ minimizing the risk $R(\theta) = \mathbb{E}_{\vec{X}, \vec{Y}} \left[\mathcal{L}_{f_{\theta}} \left(\vec{X}, \vec{Y} \right) \right]$.

Empirical risk minimization

Risk estimation

In practice, we can only estimate the risk $R(\theta)$ from a sample of (\vec{X}, \vec{Y}) .

Definition

The empirical risk for a sample $s_n = \{(\vec{x}_1, \vec{y}_1), \dots, (\vec{x}_n, \vec{y}_n)\}$ is the average loss over s_n :

$$\widehat{R}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_{f_\theta}(\vec{x}_i, \vec{y}_i)$$

Selection rule

Given a sample s_n of (\vec{X}, \vec{Y}) , select a model f_{θ} minimizing the empirical risk \hat{R}_n .

DINFK

Risk minimization flow

- 1. Select loss
- 2. Select model space
- 3. Minimize empirical risk

Loss: Prediction form

Form of the loss

For prediction it is standard to derive the loss $\mathcal{L}_{f_{\theta}}$ from a distance *d*:

 $\mathcal{L}_{f_\theta}(\vec{x},\vec{y})=d(\vec{y},f_\theta(\vec{x})).$

Shape of d

The shape of the distance d determines how difficult is to find a minimizer of \widehat{R}_n :

- The simple 0–1 distance $d(\vec{y}, \vec{y}') = (0 \text{ if } \vec{y} = \vec{y}' \text{ else } 1)$ is too difficult in practice.
- The Euclidean distance is easier since it is smooth; often suitable for regression.

Loss: Smoothing for classification

Predict k classes $\{1, ..., k\}$

- 1. The 0-1 distance most natural but difficult.
- 2. The Euclidean distance d_2 easier but biased: $y = 1 \implies d_2(y, 2) < d_2(y, 8)$.

Embed $\{1, ..., k\}$ into \mathbb{R}^k

...

$$\begin{array}{rrrr} 1 & \mapsto & (1,0,\ldots,0) \\ 2 & \mapsto & (0,1,\ldots,0) \end{array}$$

 $k \quad \mapsto \quad (0,0,\ldots,1)$

Predict probability vectors in \mathbb{R}^k

$$\mathcal{L}_{f_{\theta}}(\vec{x}, y) = d(\text{embed}(y), f_{\theta}(\vec{x}))$$

Use smooth distance d such as

1. Cross-entropy 2. KL divergence

Loss: Classification details

- 1. Pass prediction through softmax, f_{θ} = softmax $\circ g_{\theta}$, to create probability vectors.
- 2. The cross-entropy $H(\vec{p}, \vec{q})$ measures a distance from a true \vec{p} to an estimate \vec{q} .

Softmax :
$$\mathbb{R}^k \to \mathbb{R}^k$$
Cross-entropy : $\mathbb{R}^k \times \mathbb{R}^k \to \mathbb{R}$ Maps any vector to a probability vector.
softmax $(\vec{x}) = \frac{1}{e^{x_1} + \dots + e^{x_k}}(e^{x_1}, \dots, e^{x_k}).$ Inputs must be probability vectors.
 $H(\vec{p}, \vec{q}) = -\sum_{i=1}^k p_i \log(q_i).$

Full loss

$$\mathcal{L}_{f_{\theta}}(\vec{x}, \vec{y}) = H(\text{embed}(y), \text{softmax} \circ g_{\theta}(\vec{x})) \qquad f_{\theta} = \text{softmax} \circ g_{\theta}.$$

Model: Feedforward networks



Neuron: $\eta_{\vec{w},b}(\vec{x}) = \sigma(b + \vec{w} \cdot \vec{x})$

- \vec{w} : weights, *b*: bias
- σ : non-linearity

Layer: $\ell_{W,\vec{b}}(\vec{x}) = \sigma(\vec{b} + W \cdot \vec{x})$

• \vec{b} : bias vector, W: weights matrix

Network: $g_{\theta}(\vec{x}) = \ell_{\theta_m} \circ \cdots \circ \ell_{\theta_1}(\vec{x})$

•
$$\theta = [\theta_1, \dots, \theta_m]$$
: model parameters

Model: Architecture

Choices

1. Connections:

fully connected, convolutional, random

2. Non-linearities: rectifier, sigmoid, tanh, ...



steep, simple







flat for small/large inputs

18



Objective

- 1. Input: a sample of (\vec{X}, \vec{Y})
- 2. Goal: minimize $\widehat{R} = \widehat{R}_n$

Gradient descent

- 1. Select: θ_0
- 2. Iterate: $\theta_{t+1} = \theta_t \alpha_t \nabla \widehat{R}(\theta_t)$.

Minimization: Gradients



Gradient operator ∇

- 1. Input: $S : \mathbb{R}^d \to \mathbb{R}$
- 2. Output: $\nabla S : \mathbb{R}^d \to \mathbb{R}^d$ $\nabla S(\vec{u}) = \left(\frac{\partial S}{\partial u_1}(\vec{u}), \dots, \frac{\partial S}{\partial u_d}(\vec{u})\right)$

Evaluating $\nabla \widehat{R}(\vec{x})$

- 1. Backpropagation algorithm.
- 2. Fully automated in software.

Minimization: Stochastic gradient descent

Issues with gradient descent

- Gets stuck easily.
- Slow for large samples.

Variants and improvements

Momentum, Nesterov, Adagrad, AdaDetla, RMSProp, Adam, Natural gradient, ...

Stochastic gradient descent

- Subsample the sample into batches.
- Evaluate $\nabla \hat{R}$ on one batch each step.