

Introduction to
Automatic Adjoint Differentiation
- AAD -
in Machine Learning
and Finance

Antoine Savine

Modern Computational Finance

<http://antoinesavine.com>

Find the slides and examples here:

go to <http://github.com/asavine>

click 'CompFinance'

download slides 'Intro2AADinMachineLearningAndFinance.pdf'

download examples in the 'Workshop' folder

see tutorial in the 'xlCpp' folder for exporting C++ code to Excel

see basic AAD code in C++ in the file 'toyCode.h'

the rest of the repo is companion code for the book [Modern Computational Finance](#)

follow asavine on GitHub and watch the CompFinance repo to be notified of updates

Introduction

Adjoint Differentiation - AD -

- Algorithm to compute **all** differentials of a scalar function **quickly** and **accurately**

- **Quickly** (*constant time*):
Compute **all** differentials in a time similar to **one** function evaluation
- **Accurately**:
Analytic differentiation,
accurate to machine precision

$$f: \mathbb{R}^D \rightarrow \mathbb{R}$$
$$x = \begin{pmatrix} x_1 \\ \dots \\ x_D \end{pmatrix} \rightarrow y = f(x)$$

$$AD: compute: \frac{\partial f}{\partial x} = \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_D} \right) \text{ in } O(1)$$

- Application:
whenever we need to compute **many** derivatives of **one** result
- AD is *only* about speed yet a massive difference in Machine Learning and Finance – without AD:
 - Deep ANNs could not learn their parameters in reasonable time
 - Financial risks could not be computed in real time
- See article on QuantMinds:
<https://knect365.com/quantminds/article/2d9e3329-cae1-4320-86a2-a4e2303b1bfe/a-brief-introduction-to-automatic-adjoint-differentiation-aad>

Application: model fitting

- Model: $\hat{y} = f(x; \mathcal{G})$
 - Makes a prediction \hat{y} from a set of inputs $x = (x_1, \dots, x_n)^T$
 - Given parameters $\mathcal{G} = (\mathcal{G}_1, \dots, \mathcal{G}_D)^T$
- Learns its parameters by minimizing a *cost* function
 - E.g. sum of squared errors: $C(\mathcal{G}) = \sum_{i=1}^m \left[f(x^{(i)}; \mathcal{G}) - y^{(i)} \right]^2 = \sum_{i=1}^m c_i(\mathcal{G})$ (c_i : error or *loss* on training example i)
 - Over a training set of m labelled examples: $(x^{(i)}, y^{(i)})_{1 \leq i \leq m}$ ($x^{(i)} \in \mathbb{R}^n, y^{(i)} \in \mathbb{R}$)
- Must compute gradient of the cost function for optimization algorithm: $\frac{\partial C}{\partial \mathcal{G}} = \left(\frac{\partial C}{\partial \mathcal{G}_1}, \dots, \frac{\partial C}{\partial \mathcal{G}_D} \right)$
 - Standard case (sum of squared errors): $\frac{\partial C}{\partial \mathcal{G}} = \sum_{i=1}^m \frac{\partial c_i}{\partial \mathcal{G}}$, $\frac{\partial c_i}{\partial \mathcal{G}} = 2c_i \frac{\partial f(x_i; \mathcal{G})}{\partial \mathcal{G}}$

Training a deep learning model

- Multi-Layer Perceptron (MLP)
simplest and most common deep learning model:

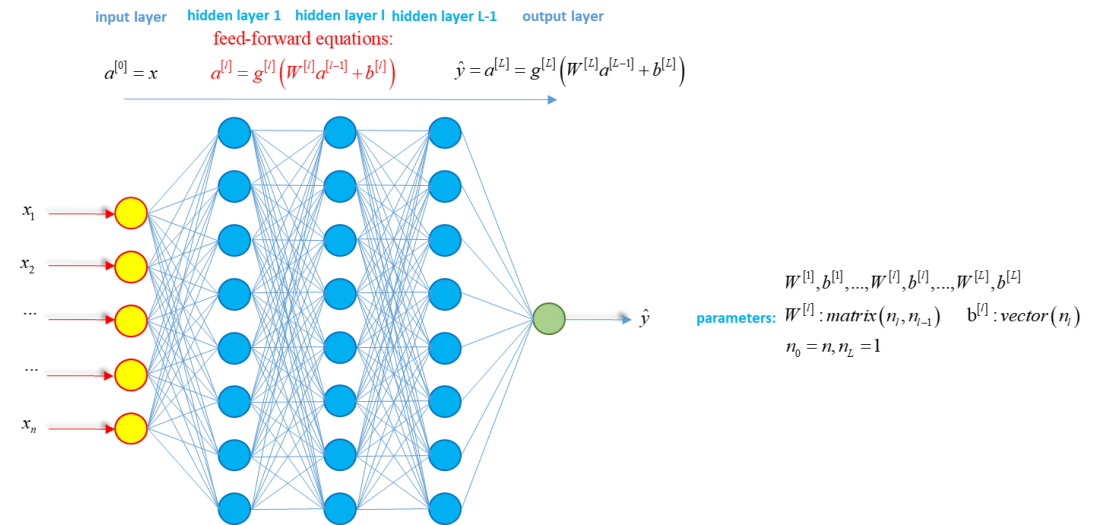
- Prediction: feed-forward equations $a^{[l]} = g^{[l]}(W^{[l]}a^{[l-1]} + b^{[l]})$
- Training set: m inputs $x^{(i)}$, each a vector in dimension n_0 with corresponding (scalar) labels $y^{(i)}$

- Train the network: learn (fit) parameters $W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}$ to minimise cost = sum of squared errors (losses) on the training set $c^{(i)} = (\hat{y}^{(i)} - y^{(i)})^2$
- Then apply trained model to make predictions on new examples with unknown label

- Requires an iterative minimization algorithm

- On each iteration, we must provide C : m predictions through a deep network

- And all their derivatives to all the parameters $\frac{\partial c^{(i)}}{\partial W_{j,k}^{[l]}}, \frac{\partial c^{(i)}}{\partial b_j^{[l]}}$



Calibrating a financial pricing model

- Calibration of a pricing model to market data:
 - Prediction: valuation of a financial product $\hat{y} = \text{product value} = f(\text{product parameters } x; \text{model parameters } \theta)$
 - Training set: m products, each with n known features $x^{(i)}$ and market price $y^{(i)}$
 - Calibration: fit parameters $\theta = (\theta_1, \dots, \theta_D)^T$ (e.g. volatilities) to minimise cost = sum of squared errors (losses) on the training set
 - Then apply calibrated model to price other products which values are unknown
- Minimization algorithms (gradient descent, conjugate gradients, pseudo-Newton, Levenberg-Marquardt...)
 - Requires derivatives $\frac{\partial C}{\partial \theta} = \left(\frac{\partial C}{\partial \theta_1}, \dots, \frac{\partial C}{\partial \theta_D} \right)$ on each iteration, we have as many differentials as we have parameters
 - in addition C may be expensive to evaluate, must value all n products in the calibration set

Notes on calibration

- Calibration is considered best practice in financial derivatives (as opposed to estimation of parameters)
- To understand exactly why, see
<https://www.slideshare.net/AntoineSavine/60-years-birthday-30-years-of-ground-breaking-innovation-a-tribute-to-bruno-dupire-by-antoine-savine>
- Calibration is only similar to training ML models in appearance
 - ML models are *entirely* trained on data and make predictions on new examples *drawn from the same distribution*
 - “Interpolation” problem
 - Financial models have a strong structure: arbitrage-free, realistic dynamics, some parameters (e.g. correlation) typically estimated
 - Often calibrated on European options to value more complicated *exotic* options
 - “Extrapolation” problem

Application: model risk

- Financial model: $value = f(\text{product features } x; \text{model parameters } \mathcal{G})$
- Model parameters:
 - Today's underlying asset prices
 - Volatilities, correlations
 - Mean-reversions
 - Etc.
- Model risks: $\frac{\partial f}{\partial \mathcal{G}}$ with potentially massive number of parameters
 - Black-Scholes: spot, volatility, rate
 - Dupire: two-dimensional local volatility surface, with stochastic volatility: vol of vol, correlation spot/vol
 - Interest Rate Models: initial yield curve + spread curve(s) + volatility curve or surface (+ mean-reversion...)
 - Multi-currency models (Baskets, CVA...): all of the above for all currencies + forex rates and their volatilities + giant correlation matrix
 - Typically in the thousands or tens of thousands
 - Valuation generally expensive: Monte-Carlo simulations

Market risk

- Usually, we are not interested in risks to model parameters but risks to market variables

E.g. Dupire: not local volatilities but implied volatilities

Model parameters are generally calibrated to market variables

a : vector of k market variables (underlying asset prices, rates, spreads, implied volatilities...)

b : vector of n model parameters (calibrated to market variables so $b = c(a), c: \mathbb{R}^k \rightarrow \mathbb{R}^n$)

g : valuation function of model (generally numerical such as Monte-Carlo simulation) $V = g(b), g: \mathbb{R}^n \rightarrow \mathbb{R}$

f : valuation function of market $V = g(b) = g[c(a)] = f(a), f: \mathbb{R}^k \rightarrow \mathbb{R}$

- Risk report computes model risks $\frac{\partial g}{\partial b}$ and/or market risks $\frac{\partial f}{\partial a}$

- Market risks computed from model risks: $a \in \mathbb{R}^k \xrightarrow{c} b \in \mathbb{R}^n \xrightarrow{g} v \in \mathbb{R}$ so by the chain rule $\frac{\partial f}{\partial a} = \frac{\partial g}{\partial b} \frac{\partial c}{\partial a}$

- To compute $\frac{\partial c}{\partial a}$ and market risks $\frac{\partial f}{\partial a} = \frac{\partial g}{\partial b} \frac{\partial c}{\partial a}$ knowing model risks $\frac{\partial g}{\partial b}$, see for instance:

http://papers.ssrn.com/sol3/papers.cfm?abstract_id=3262571

- In this presentation, we focus on model risks $\frac{\partial g}{\partial b}$

Differentiation

- In all these applications (and many others)
 - We have a scalar function f of many inputs $f: \mathbb{R}^D \rightarrow \mathbb{R}$
 - This function is typically expensive to evaluate
Losses of deep neural nets on large training set or a Monte-Carlo valuation may take several seconds, even on parallel hardware
 - It typically takes a large number D of inputs
Thousands of weights for deep neural nets, or thousands of market/model variables for financial valuation
 - We must compute its sensitivities to all its inputs $\frac{\partial f}{\partial x} = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_D} \right)$
- Conventional differentiation is of linear complexity in the number of inputs/differentials
 - To compute **one** differential takes time similar to **one** evaluation of the function
 - For example, with (one-sided) finite differences $\frac{\partial f}{\partial x_i} \approx \frac{f(x + \varepsilon e_i) - f(x)}{\varepsilon}, e_i = (\delta_{ij})_{1 \leq j \leq n}$
 - We have $(D+1)$ function evaluations to compute D differentials
 - Illustration: if one function evaluation takes 1sec and we have 1,000 inputs
It takes 1,001sec (15+ minutes) to evaluate its gradient once
 - This is not viable, a different technology is necessary

Adjoint Differentiation - AD -

- Algorithm to compute analytically all the differentials of f : $\frac{\partial f}{\partial x} = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_D} \right)$
 - In a time independent of D
 - With an efficient implementation, it takes 4 to 10 times **one** evaluation of f to compute **all** its differentials
 - So AD can compute the 1,000 differentials of a function that takes 1sec in ~4 to 10sec as opposed to 15min+
 - Without loss of accuracy (if anything, AD differentials are *more* accurate)
- Hence the importance of this technology today
 - In finance, AD computes thousands of complex risks accurately and very quickly
 - In machine learning, AD computes the gradient of cost functions to many parameters in constant time
 - The same technology powers banks to estimate financial risks and deep nets to learn their weights
 - AD
 - Allows Banks to compute massive amount of regulatory risks in real time
 - Powers your telephones to learn to recognize you and your friends in reasonable time
 - Is largely credited for the recent and spectacular successes in the field of deep learning

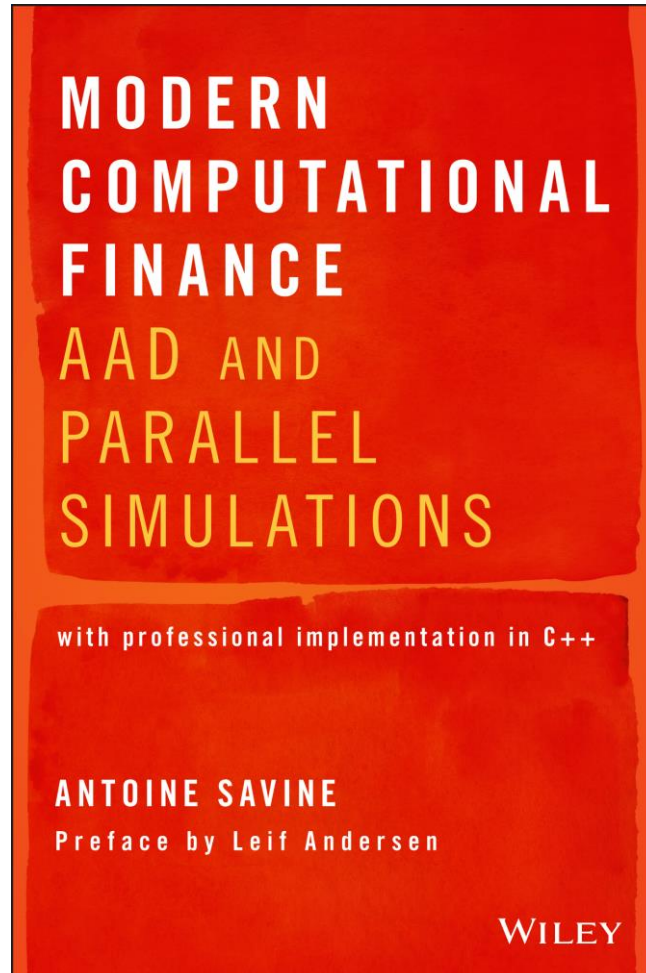
A brief history of AD

- Invented in the 1960s (Wengert, 1964)
- Called “holy grail” of sensitivity computation (Griewank, 2012)
- Classified as one of the 30 greatest numerical algorithms of the 20th century (Trefethen, 2015)
- Did not take firm hold in the computer science community until the 1980s
- Applied in deep learning to efficiently compute gradients of cost functions in learning algorithms
- Known under the names
Back-Propagation (or simply back-prop, in deep learning),
reverse differentiation, backward differentiation, adjoint accumulation, algorithmic differentiation, etc.
- Not adopted in Finance until 2006 (Giles and Glasserman’s “Smoking Adjoint”)
- Large scale implementation in Danske Bank with parallel Monte-Carlo simulations won In-House System of the Year 2015 Risk Award
- Widely adopted today for risk and calibration
- Delays in adoption mainly due to complexity, lack of teaching material and challenges in practical implementation (Andersen, 2018)

Automatic Adjoint Differentiation - AAD -

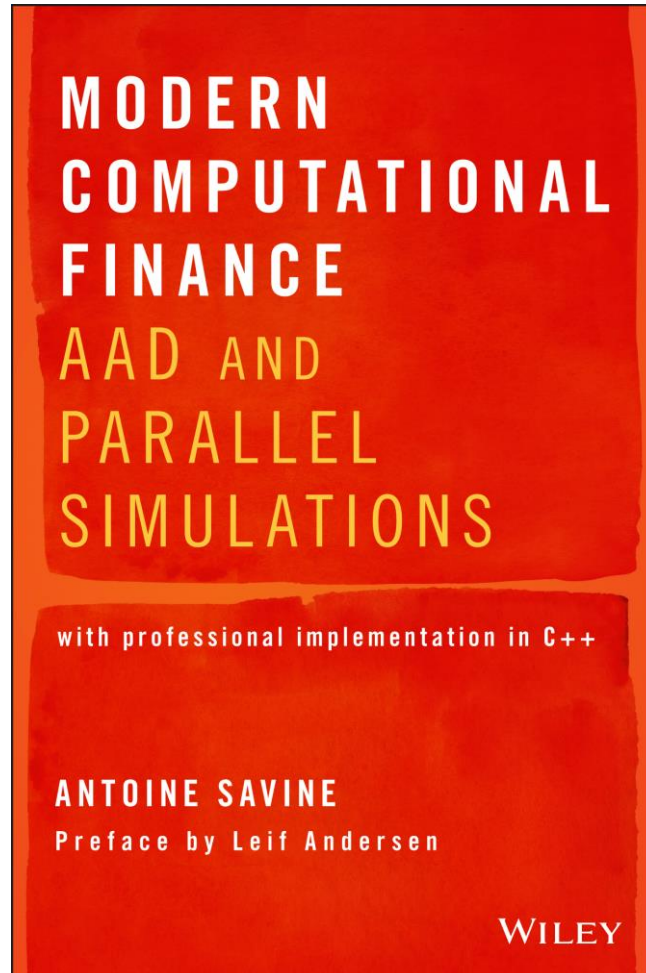
- Automatic implementation of AD
 - Developers only produce the evaluation code for the function f
 - The framework automatically produces (constant time) differentiation “code” for $\partial f / \partial x$
 - So developers don’t need to write (complicated and error prone) AD code
 - And differentiation code is automatically updated when evaluation code changes
- Many commercial and open source AAD frameworks exist
 - Some are generic, some specialize in finance or machine learning
 - Some are free, some as extremely expensive
 - Work in a variety of programming languages: Python, C++, ...
 - Work either with source transformation, operator overloading or by letting clients build evaluation graphs
- Example: TensorFlow, one of most popular frameworks in machine learning
 - Written in C++/CUDA (multicore CPU/GPU)
 - With APIs in Python, JavaScript, Java, Go, Swift...
 - Allow clients to create evaluation graphs
 - TensorFlow evaluates the graphs and automatically applies AD to compute derivatives through the graphs

AAD book



- Written by (some of) the people who wrote Danske Bank's award winning systems
- Prefaced by Leif Andersen, read preface on researchgate.net/publication/328042479_Modern_Computational_Finance_AAD_and_Parallel_Simulations
- Teaches:
 - The modern design and efficient implementation of financial simulation libraries
 - The implementation of *parallel* simulation libraries
 - And of course AAD
- AAD is covered in deep detail:
 - Conceptual and mathematical foundations
 - Practical implementation
 - Memory management and check-pointing
 - Application in the context of large Monte-Carlo simulations, including in parallel
 - Application to model and market risks
 - Cutting-edge implementation with meta-programming and expression templates
 - And much more
- Ships with complete, professional code in C++
 - Explained in deep detail in the book
 - Freely available on the GitHub repo <http://www.github.com/asavine/CompFinance/wiki>

AAD book reviews



From Amazon: <https://www.amazon.com/gp/product/1119539455>

It would not be much of an exaggeration to say that Antoine Savine's book ranks as the 21st century peer to Merton's 'Continuous-Time Finance'.

Vladimir Piterbarg

This book [...] addresses the challenges of AAD head on. [...] The exposition is [...] ideal for a Finance audience. The conceptual, mathematical, and computational ideas behind AAD are patiently developed in a step-by-step manner, where the many brain-twisting aspects of AAD are demystified. For real-life application projects, the book is loaded with modern C++ code and battle-tested advice on how to get AAD to run for real. [...] Start reading!

Leif Andersen

An indispensable resource for any quant. Written by experts in the field and filled with practical examples and industry insights that are hard to find elsewhere, the book sets a new standard for computational finance.

Paul Glasserman

Overview

AAD: Demonstration

1. Adjoint Differentiation for Deep Learning
 1. A brief introduction to Artificial Neural Networks (ANN)
 2. Back-Propagation through ANNs
2. Adjoint Differentiation for arbitrary calculations
 1. Evaluation graphs
 2. Back-Propagation through evaluation graphs
3. Automatic Adjoint Differentiation for financial derivatives in (simple) C++
 1. AAD with operator overloading
 2. AAD over Monte-Carlo simulations

Demonstration: Dupire's model (1992)

- Extended Black & Scholes dynamics (in the absence of rates, dividends etc.): $\frac{dS}{S} = \sigma(S, t) dW$
- Calibrated with Dupire's celebrated formula: $\sigma(K, T) = \frac{2 \frac{\partial C}{\partial T}}{\frac{\partial^2 C}{\partial K^2}}$ with $C(K, T)$ = call prices of strike K , maturity T
- Implemented with a (bi-linearly interpolated) local volatility matrix: $\sigma_{ij} = \sigma(S_i, T_j)$
- Volatility matrix: 21 spots (every 5 points 50 to 200) and 36 times (every month from now to 3y)
we have 1,080 volatilities + 1 initial spot (=100) = 1,081 model parameters
- Valuation of a 3y (weekly monitored) barrier option strike $K=120$, barrier $B=150$: $v(S_t, t) = E \left[(S_T - K)^+ 1_{\{\max(S_{T_1}, \dots, S_{T_K}) < B\}} \middle| S_t \right]$
- Solved with Monte-Carlo or FDM over the equivalent PDE (from Feynman-Kac's theorem)
- We focus on Monte-Carlo simulations here: 500,000 paths, 156 (weekly) time steps

Demonstration: Results

- Implementation
 - C++ code exported to Excel (see [tutorial ExportingCpp2xl.pdf](#) on the GitHub repo, folder `xlCpp`)
 - Generic library design with efficient implementation (Chapter 6)
 - Parallel implementation (Chapters 3 and 7)
 - Sobol quasi-random numbers (Chapters 5 and 6), excerpt here: https://medium.com/@antoine_savine/sobol-sequence-explained-188f422b246b
 - Advanced AAD with expression templates (Chapter 15)
- Hardware: quad-core laptop (surface book 2, 2017)
- Performance:
 - One evaluation with 500,000 paths over 156 time steps take ~0.8sec
 - We have 1,081 risk sensitivities, take about 15 minutes to produces model risk report with linear differentiation
 - With AAD the 1,081 differentials are produced in ~1.5sec

Demonstration material

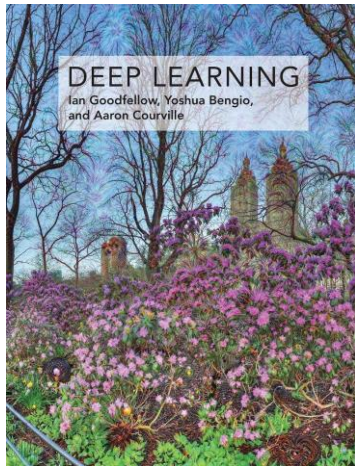
- GitHub repo (<http://github.com/asavine/CompFinance>)
- Prebuilt Xl addin **xlComp.xll** (may need to install redistributables VC_redist.x86.exe and VC_redist.x64.exe)
- Demonstration spreadsheet **xlTest.xlsx**

[illegible][illegible]

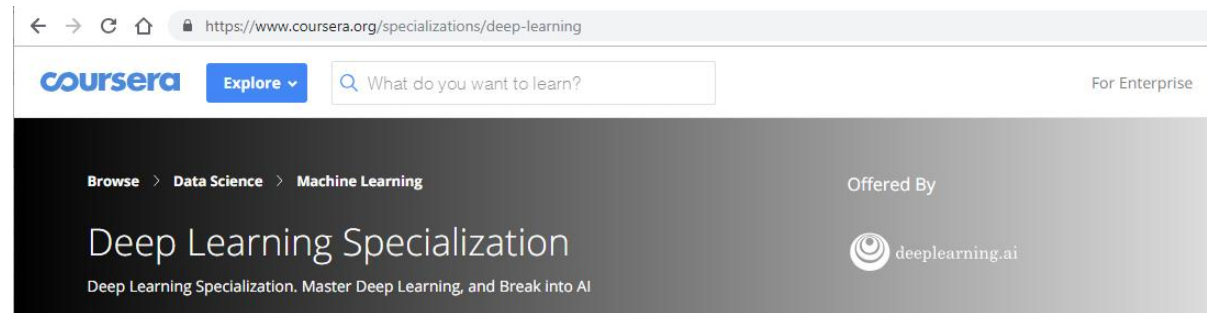
Deep Learning

Neural networks and deep learning

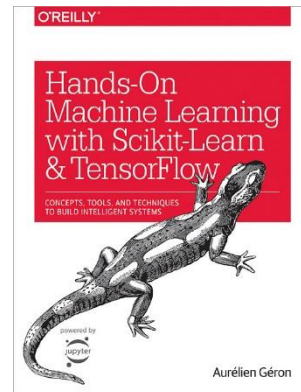
- Adjoint differentiation is best known and explained in the context of deep learning
- We therefore *briefly* introduce deep learning
- (Much) deeper learning material is found, for example, in:



Goodfellow's reference book



Andrew Ng's videos on Coursera



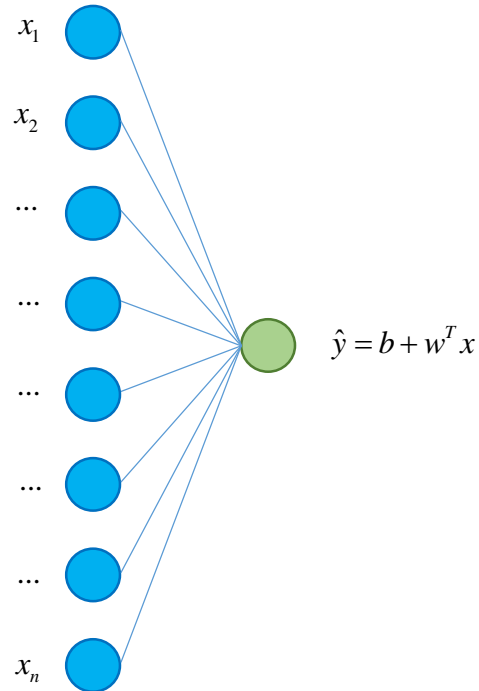
Geron's practical introduction with TensorFlow

...and many other books and resources

Linear regression

- Linear model (joint Gaussian assumptions): $\hat{y} = E[y|x] = b + \sum_{i=1}^n w_i x_i = b + w^T x$

- Parameters: $b \in \mathbb{R}$ and $w \in \mathbb{R}^n$

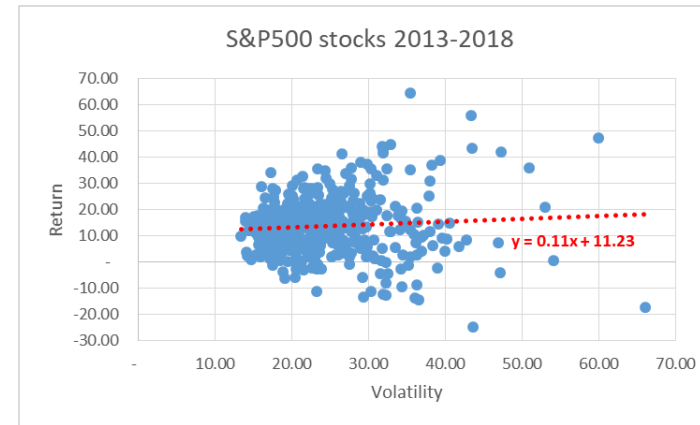


Linear regression: example

- In *regression* problems we predict real numbers

Example: predict return from volatility

- Each point is a stock in the S&P500
- Horizontal axis is annual volatility
- Vertical axis is annual return
- Estimated with daily data, 2013-2018
- Evidence of (small) risk premium
- Poor regression quality



Linear classification

- Alternatively, *classification* problems predict discrete categories

Example: identify animals in pictures: 0: not an animal, 1: cat, 2: dog, 3: bird, 4: other animal

picture 1920x1080x24bit colour



vector of 1920x1080 pixels

$$\in [0, 2^{24} - 1]$$

$$x = \begin{pmatrix} x_1 \\ \dots \\ x_{2,073,600} \end{pmatrix}$$

softmax regression

$$\hat{y} = s(b + wx) = \begin{bmatrix} \Pr(\text{no animal}) \\ \Pr(\text{cat}) \\ \Pr(\text{dog}) \\ \Pr(\text{bird}) \\ \Pr(\text{other animal}) \end{bmatrix}$$

Parameters:

b: vector of dimension 5

w: matrix 5 x 2,073,600

s: softmax function

$s: \mathbb{R}^5 \rightarrow (0,1)^5$ summing to 1

$$s(z) = \begin{bmatrix} \frac{e^{z_i}}{\sum e^{z_j}} \end{bmatrix}$$

- Regression more relevant than classification in quantitative finance
- Deep learning community generally focuses more on classification
- See literature for details, for example Stanford's CS229 on <http://cs229.stanford.edu/syllabus.html>
- This presentation sticks with regression
- Everything that follows generalises easily to classification and the maths remain essentially identical

Linear fit

- Training set of m examples $x^{(1)}, \dots, x^{(m)}$ (each a vector in dimension n) with corresponding labels $y^{(1)}, \dots, y^{(m)} \in \mathbb{R}$

- Learn parameters: $b \in \mathbb{R}$ and $w \in \mathbb{R}^n$ by minimizing prediction errors (cost function) $C(b, w) = \sum_{i=1}^m \left(\underbrace{b + w^T x^{(i)}}_{=\hat{y}^{(i)}} - y^{(i)} \right)^2$

- Note that this is the same as maximizing (log) likelihood under Gaussian assumptions, hence:

$b^*, w^* = \arg \min C(b, w)$ are the maximum likelihood estimators (MLE) of the parameters

- b^* and w^* are found analytically, solving for $\frac{\partial C}{\partial b} = 0$ and $\frac{\partial C}{\partial w} = 0$

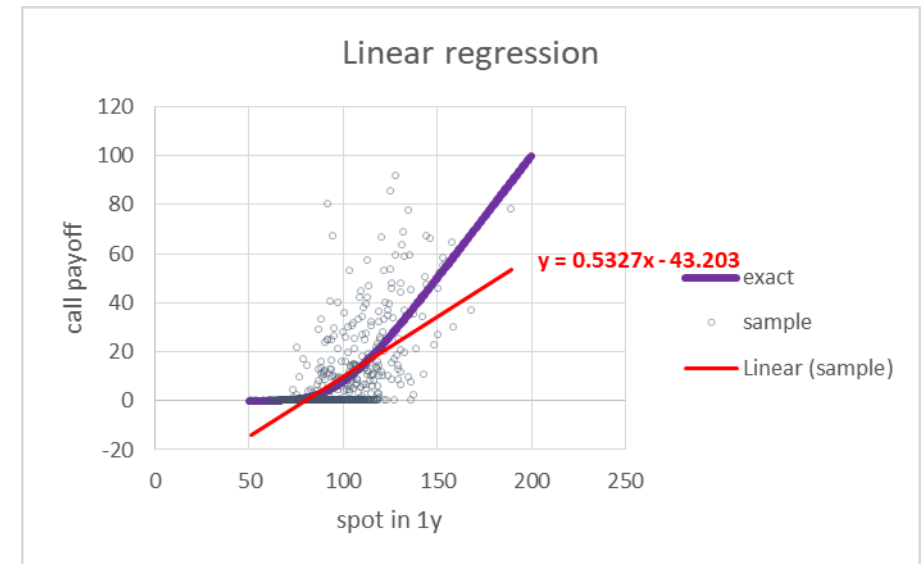
- Result (“normal equation”): $\begin{pmatrix} b \\ w_1 \\ w_2 \\ \dots \\ w_n \end{pmatrix} = (X^T X)^{-1} X^T Y$ where $X = \begin{pmatrix} 1 & x_1^{(1)} & \dots & x_n^{(1)} \\ 1 & x_1^{(2)} & \dots & x_n^{(2)} \\ \dots & x_1^{(i)} & x_j^{(i)} & x_n^{(i)} \\ 1 & x_1^{(m)} & \dots & x_n^{(m)} \end{pmatrix}$ and $Y = \begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \dots \\ y^{(m)} \end{pmatrix}$

Linear limitations

- Important revaluation example:

- Predict the future price in 1y of a European call strike 100, maturity 2y
- By regression of the payoff in 2y: $y^{(i)} = (S_{2y}^{(i)} - K)^+$
- Over the underlying asset price in 1y: $x^{(i)} = S_{1y}^{(i)}$
- With a training set of m paths $(S_{1y}^{(i)}, S_{2y}^{(i)})_{1 \leq i \leq m}$
generated under Black & Scholes' model (spot = 100, volatility of 20%)

- We know the exact solution is given by Black and Scholes' formula so we can assess the quality of the regression

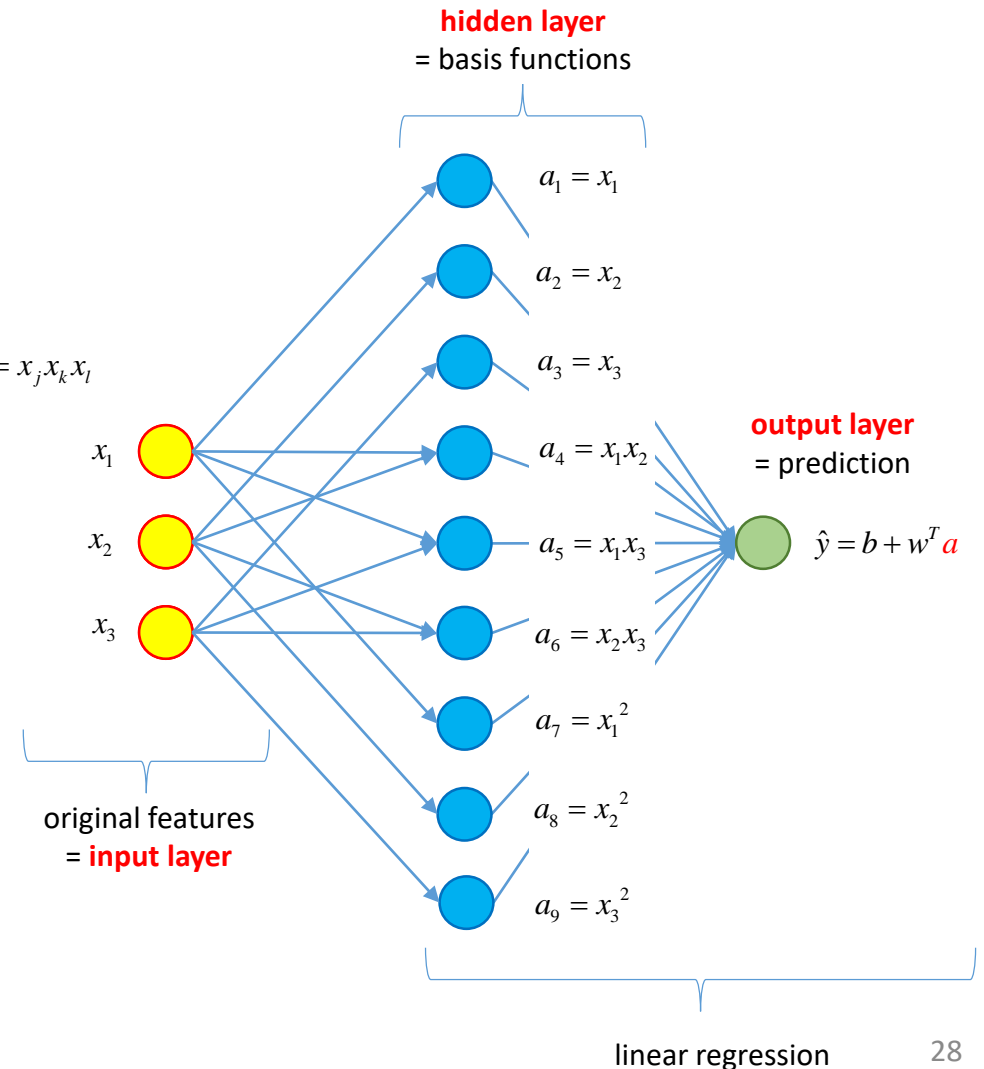


- See spreadsheet reval.xlsx on the repo <http://github.com/asavine/CompFinance/Workshop/>
- Linear regression obviously fails to approximate the correct function because it cannot capture non-linearities

Basis function regression

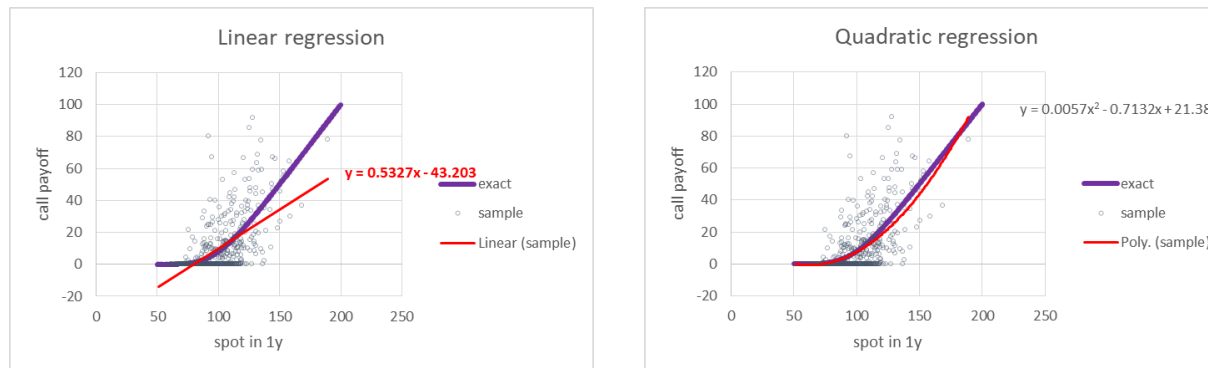
- Solution: regress not on x but on basis functions a of x
- Example: polynomial regression
 - Basis functions = monomials (x is in dimension n_0):
 - 1st degree: all the x s $a_i = x_i$
 - 2nd degree: all the squares $a_i = x_j^2$ and pair-wise products $a_i = x_i x_j$
 - 3rd degree: all the cubes $a_i = x_j^3$ and pair-wise $a_i = x_j x_k^2$ and triplet-wise $a_i = x_j x_k x_l$
 - Etc.
- Prediction in two steps:
 - Start with the vector x of n_0 features
 - Compute vector of n_1 basis functions: $a = \varphi(x)$
 - Predict linearly **in the basis functions**: $\hat{y} = b + w^T a$
- Training: identical to linear regression on a in place of x

$$\begin{pmatrix} b \\ w_1 \\ w_2 \\ \dots \\ w_{n_1} \end{pmatrix} = (A^T A)^{-1} A^T Y \quad \text{where} \quad A = \begin{pmatrix} 1 & a_1^{(1)} & \dots & a_{n_1}^{(1)} \\ 1 & a_1^{(2)} & \dots & a_{n_1}^{(2)} \\ \dots & a_1^{(i)} & a_j^{(i)} & a_{n_1}^{(i)} \\ 1 & a_1^{(m)} & \dots & a_{n_1}^{(m)} \end{pmatrix} \quad \text{and} \quad a^{(i)} = \varphi(x^{(i)})$$



Still linear, no more limited

- Works nicely: quadratic regression of $(S_{T_2} - K)^+$ over S_{T_1} and S_{T_2} approximates Black & Scholes' formula well in simulated example



- Remarkable how the algorithm manages to detect Black & Scholes' pattern in noisy data
- Mathematically:
 - Combinations of polynomials can approximate any smooth function to arbitrary precision
 - Hence, detect any (smooth) non-linear pattern in data
 - *But only with a large number of basis functions / high polynomial degree*

Curse of dimensionality

- How many monomials in a p -degree polynomial regression?

Number n_1 of basis functions (dimension of a) grows exponentially in number n_0 of features (dimension of x)

Precisely: $n_1 = \frac{(n_0 + p)!}{n_0! p!} - 1$

- Quadratic regression: $n_1 = \frac{(n_0 + 2)(n_0 + 1)}{2} - 1$ dimension grows e.g. from 10 to 65, from 100 to 5,151, from 1,000 to 501,500
- Cubic regression: $n_1 = \frac{(n_0 + 3)(n_0 + 2)(n_0 + 1)}{6} - 1$ dim grows 10 to 286, 100 to 176,851, 1,000 to 167,668,501
- In general, number of basis functions increases exponentially in dimension
- “Rule of ten”:
to avoid overfitting small data set with many parameters, we need an amount of data (training examples) of approximately $m > 10n$
- Hence, basis function regression requires amount of data exponential in (original) dimension n_0
- Basis function regression works nicely in low dimension but doesn’t scale (think of image processing where $n_0 = 1920 \times 1080 > 2M$)

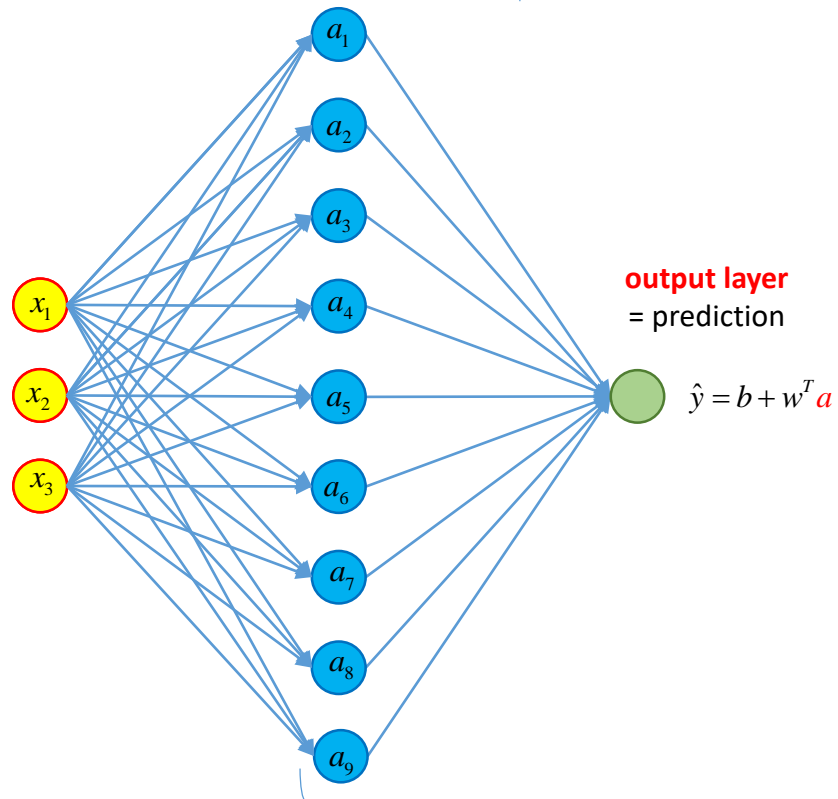
Introduction to overfitting

- If you fit a linear model with n free parameters to n data points
 - You will find a perfect fit
 - But the model may not generalize well to new examples
 - Because it captured the noise of the particular training set, “overfitting”
- For error analysis (“variance-bias trade-off”), see:
<https://www.quora.com/Does-neural-network-generalize-better-if-it-is-trained-on-a-larger-training-data-set-Are-there-any-scientific-papers-showing-it/answer/Antoine-Savine-1>
- Solutions exist such as regularization
See **Stanford’s Machine Learning class on Coursera** or **Bishop’s Pattern Recognition and Machine Learning**
- But the most effective means to avoid overfitting is increase the size of the learning set
- An empirical rule of thumb called “rule of ten”
recommends a data set of at least 10x the number of fitted parameters

ANN: learn basis functions from data

BASIS FUNCTION REGRESSION

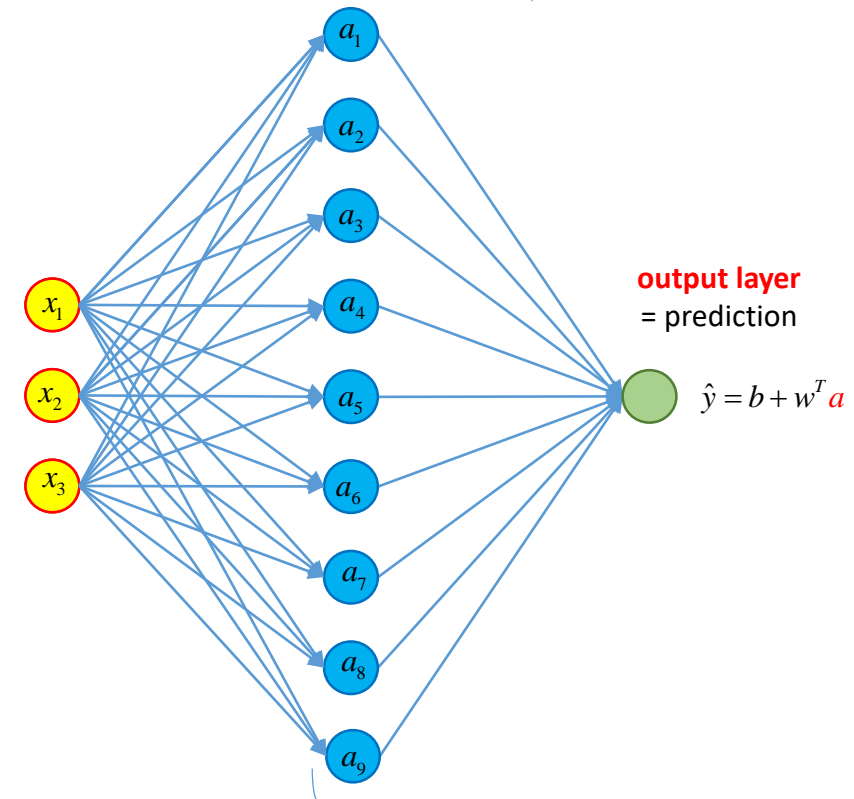
hidden layer = **arbitrary** set of basis functions
pre-processing = not part of the model, no parameters



linear regression over basis functions

ANN

hidden layer = **learned** set of basis functions
part of the model, subject to trainable parameters



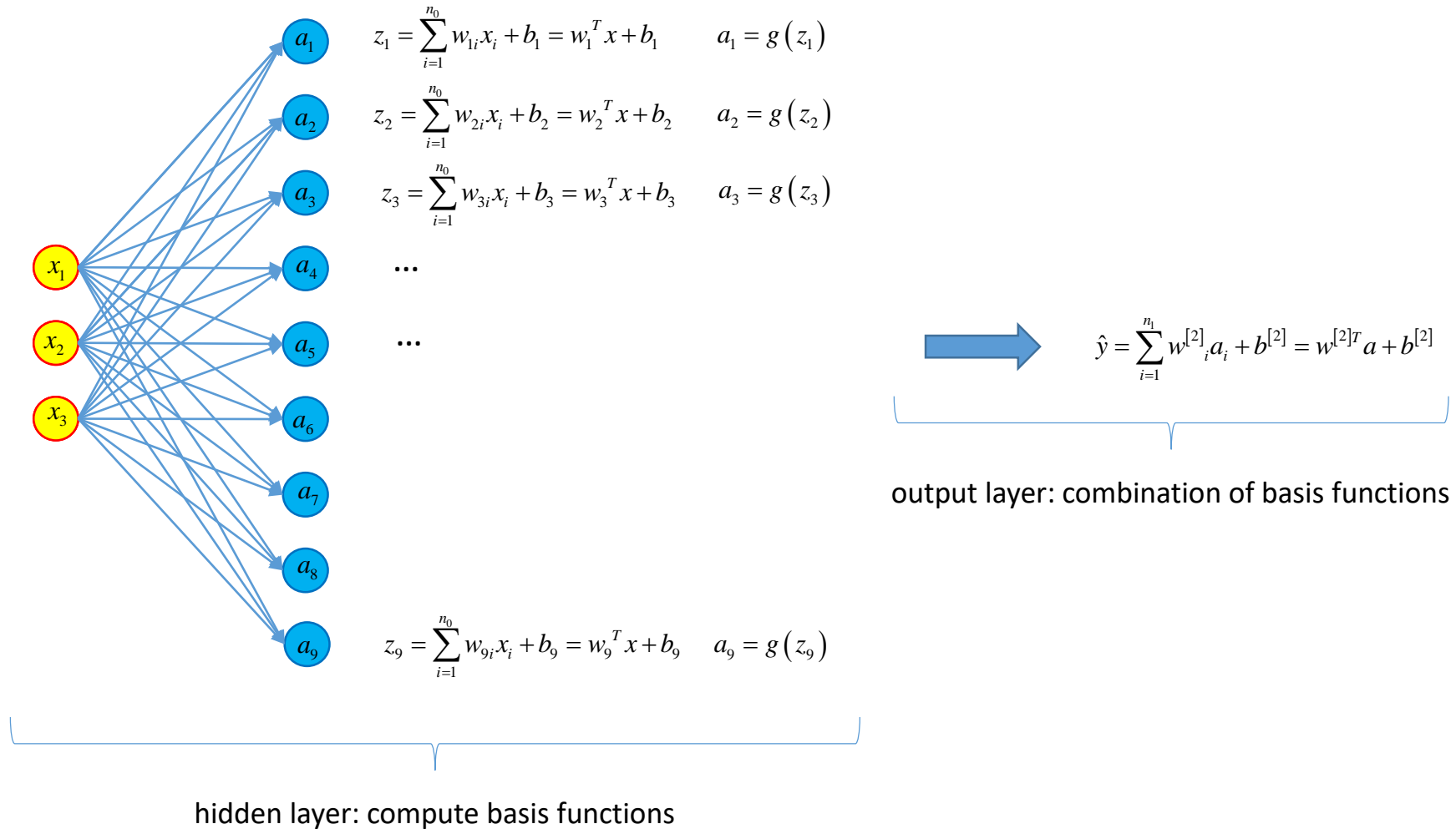
linear regression over basis functions

identical

Automatic feature extraction

- How do we learn basis functions a from the data?
 - Parametric family of basis functions: $a_i = \varphi(x; \mathcal{G}_i)$ \rightarrow each basis function is a parametric function of *the same vector of inputs x each with different parameters θ_i*
 - Learn parameters \mathcal{G} (along with the weights w and bias b of regression) by minimization of the cost (e.g. sum of squared errors)
- Most common choice of parametric basis functions: “perceptron”
 - Each basis function is a non-linear scalar function g of a linear combination of the x s (different for each basis function): $a_i = g(w_i^T x + b_i)$
 - Where each $w_i, 1 \leq i \leq n_1$ is a vector in dimension n_0 and each $b_i, 1 \leq i \leq n_1$ is a real number
 - Hence, $W = \begin{pmatrix} w_1^T \\ \dots \\ w_{n_1}^T \end{pmatrix}$ is a $n_1 \times n_0$ matrix and $b = \begin{pmatrix} b_1 \\ \dots \\ b_{n_1} \end{pmatrix}$ is a n_1 dimensional vector
 - Then the vector of basis functions is $a = g(Wx + b)$ where g is a scalar function applied element-wise to the n_1 vector $z = Wx + b$

Perceptron



Prediction with perceptrons

- Prediction in two steps

- Step 1 (hidden layer): feature (=basis function) extraction **inputs x (dim n_0) \rightarrow basis functions a (dim n_1)** $a = g(W^{[1]}x + b^{[1]})$
 - Parameters: $W^{[1]}$ matrix in dimension $n_1 \times n_0$ and $b^{[1]}$ vector in dimension n_1
 - The *activation* function g is scalar and applied element-wise on the n_1 vector $z^{[1]} = W^{[1]}x + b^{[1]}$
 - g must be non-linear or we are back to linear regression (a linear combination of linear functions is a linear function of the inputs)
- Step 2: (output layer): regression **basis functions a (dim n_1) \rightarrow prediction (real number)** $\hat{y} = w^{[2]}a + b^{[2]}$
 - Parameters: $w^{[2]}$ column vector in dimension n_1 and $b^{[2]}$ real number

- Unified notation for both steps = **feed-forward equation** $a^{[l]} = g^{[l]}(W^{[l]}a^{[l-1]} + b^{[l]})$ or $z^{[l]} = W^{[l]}a^{[l-1]} + b^{[l]} \rightarrow a^{[l]} = g^{[l]}(z^{[l]})$

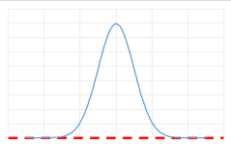
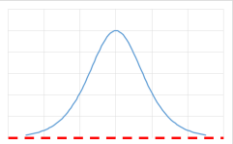
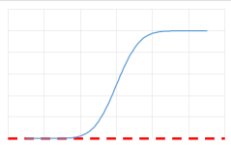
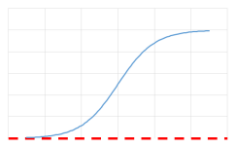
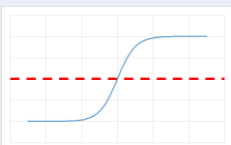
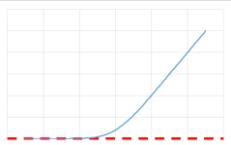
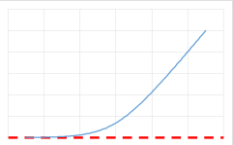
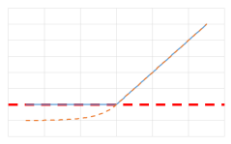
- With the notations

$a^{[0]} = x$	input layer	dim n_0	
$a^{[1]} = a$	hidden layer	dim n_1	
$\hat{y} = a^{[2]}$	output layer	dim $n_2 = 1$	
$g^{[1]} = g, g^{[2]} = id$	hidden and output activations		

 and parameters $W^{[l]}$ of dim $n_l \times n_{l-1}$ and $b^{[l]}$ of dim n_l for $l = 1, 2$

Activations

The hidden activation function g must be non-linear – here are the most common choices :

Context	Intuitive solution	Cheaper computation	Improves training*
Activate localized spheres Application: interpolation/spline	Gaussian density $g(z) = n(z) = \frac{e^{-\frac{z^2}{2}}}{\sqrt{2\pi}}$ 	Derivative of sigmoid $g(z) = \sigma'(z) = \sigma(z)[1 - \sigma(z)]$ 	
Activate above threshold Application: probability/classification/flat extrapolation	Gaussian distribution $g(z) = N(z) = \int^z n(t) dt$ 	Sigmoid $g(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$ 	Hyperbolic tangent $g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$ 
Activate proportionally to distance to threshold Application: regression with linear extrapolation	Integral of Guassian (Bachelier) $g(z) = \int^z N(t) dt$ 	SoftPlus = Integral of sigmoid $g(z) = \int^z \sigma(t) dt = \log(1 + e^z)$ 	RELU / ELU $RELU : g(z) = \max(0, z)$ $ELU : g(z) = z \text{ if } z > 0$ else $e^z - 1$ 

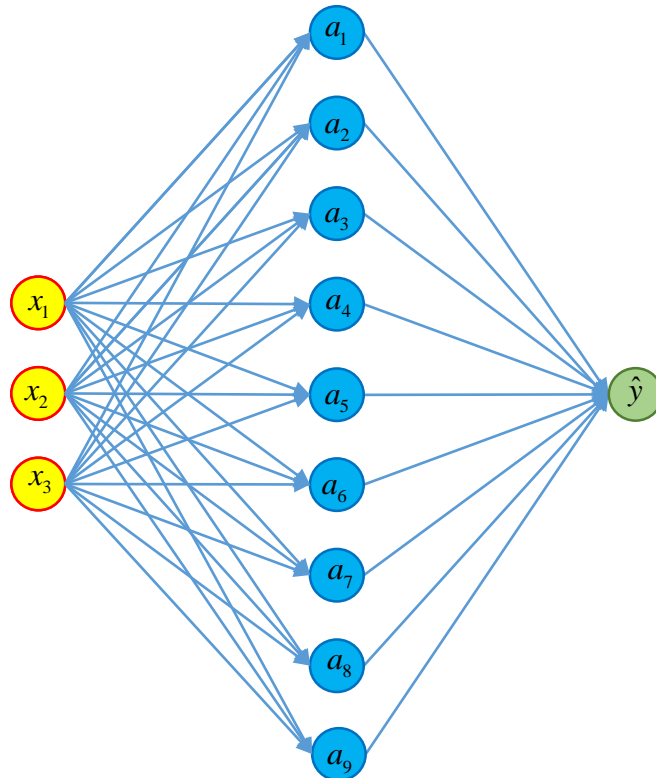
derivative
 integral

*antisymmetric around zero, saturates slowly on the sides

default choice is ELU

Computation graph

$$\begin{array}{ccc} \text{Input layer } l=0 & \text{Hidden layer } l=1 & \text{Output layer } l=2 \\ a^{[0]} = x & \longrightarrow a^{[1]} = g^{[1]}(W^{[1]}a^{[0]} + b^{[1]}) & \longrightarrow a^{[2]} = g^{[2]}(W^{[2]}a^{[1]} + b^{[2]}) = \hat{y} \\ n_0 \text{ units} & n_1 \text{ units} & n_2 = 1 \text{ units} \end{array}$$



feed-forward equation

$$a^{[l]} = g^{[l]}(W^{[l]}a^{[l-1]} + b^{[l]})$$

Computation complexity (time)

To the leading order (ignoring activations and additions)

We have 2 matrix by vector products $W^{[1]}a_0$ and $W^{[2]}a_1$
 $n_1 \times n_0 \quad n_0 \quad n_2 \times n_1 \quad n_1$

Hence quadratic complexity $n_1n_0 + n_2n_1$

To simplify when all layers have n units, complexity $\sim 2n^2$

More generally, with L hidden layers, complexity $\sim Ln^2$

Complexity is linear in layers, quadratic in units

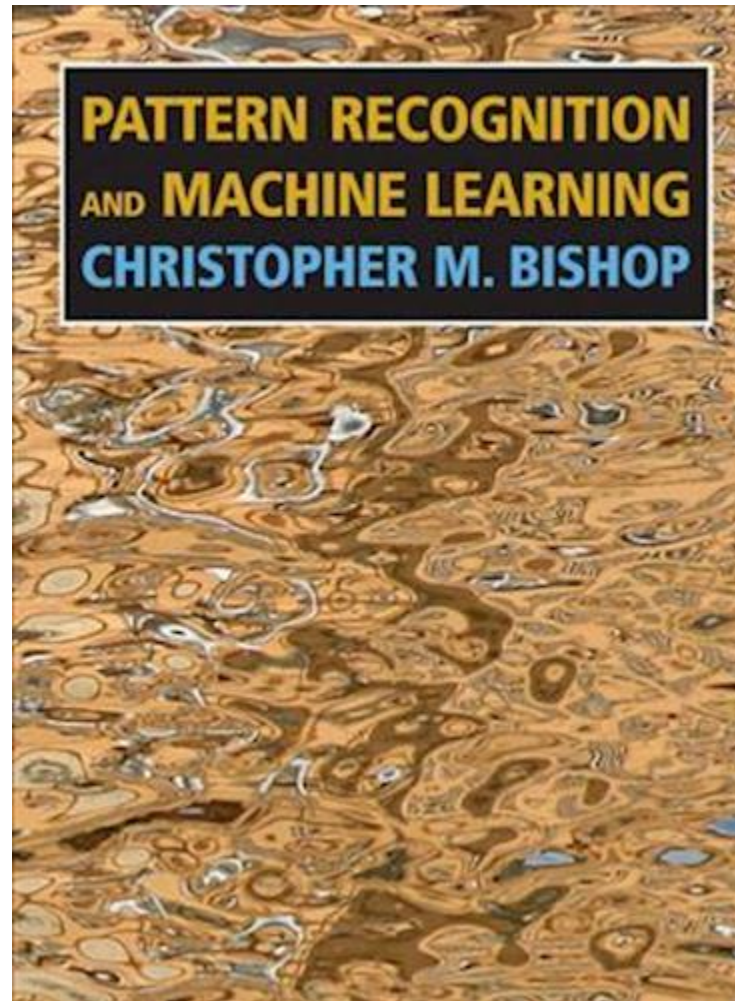
How to train your perceptron

- Training set of m examples $x^{(1)}, \dots, x^{(m)}$ (each a vector in dimension n_0) with corresponding labels $y^{(1)}, \dots, y^{(m)} \in \mathbb{R}$
- Learn parameters: $b^{[1]} \in \mathbb{R}^{n_1}, b^{[2]} \in \mathbb{R}^{n_2=1}$ and $W^{[1]} \in \mathbb{R}^{n_1 \times n_0}, W^{[2]} \in \mathbb{R}^{(n_2=1) \times n_1}$ that is $D = n_1 + n_2 + n_1 n_0 + n_2 n_1$ parameters
- By minimizing the cost function $C(b^{[1]}, b^{[2]}, W^{[1]}, W^{[2]}) = \sum_{i=1}^m c_i(b^{[1]}, b^{[2]}, W^{[1]}, W^{[2]})$ with loss $c_i(b^{[1]}, b^{[2]}, W^{[1]}, W^{[2]}) = \left(a^{[2](i)}_{=\hat{y}^{(i)}} - y^{(i)} \right)^2$
- Note we learn the basis functions and the regression weights at the same time
- No more analytic solution (C is not even convex in the D parameter space...)
- We must apply a numerical minimization algorithm
- Many iterative algorithms exist in literature and standard libraries
- The simplest one is gradient descent:
repeatedly move parameters one step down the steepest slope of the cost function (gradient of cost to parameters)
- Gradient descent and other algorithms require differentials of cost function to all the parameters
- On each iteration, we must compute not only the cost function but also its D derivatives where we recall that $D \sim O(n_1 n_0)$

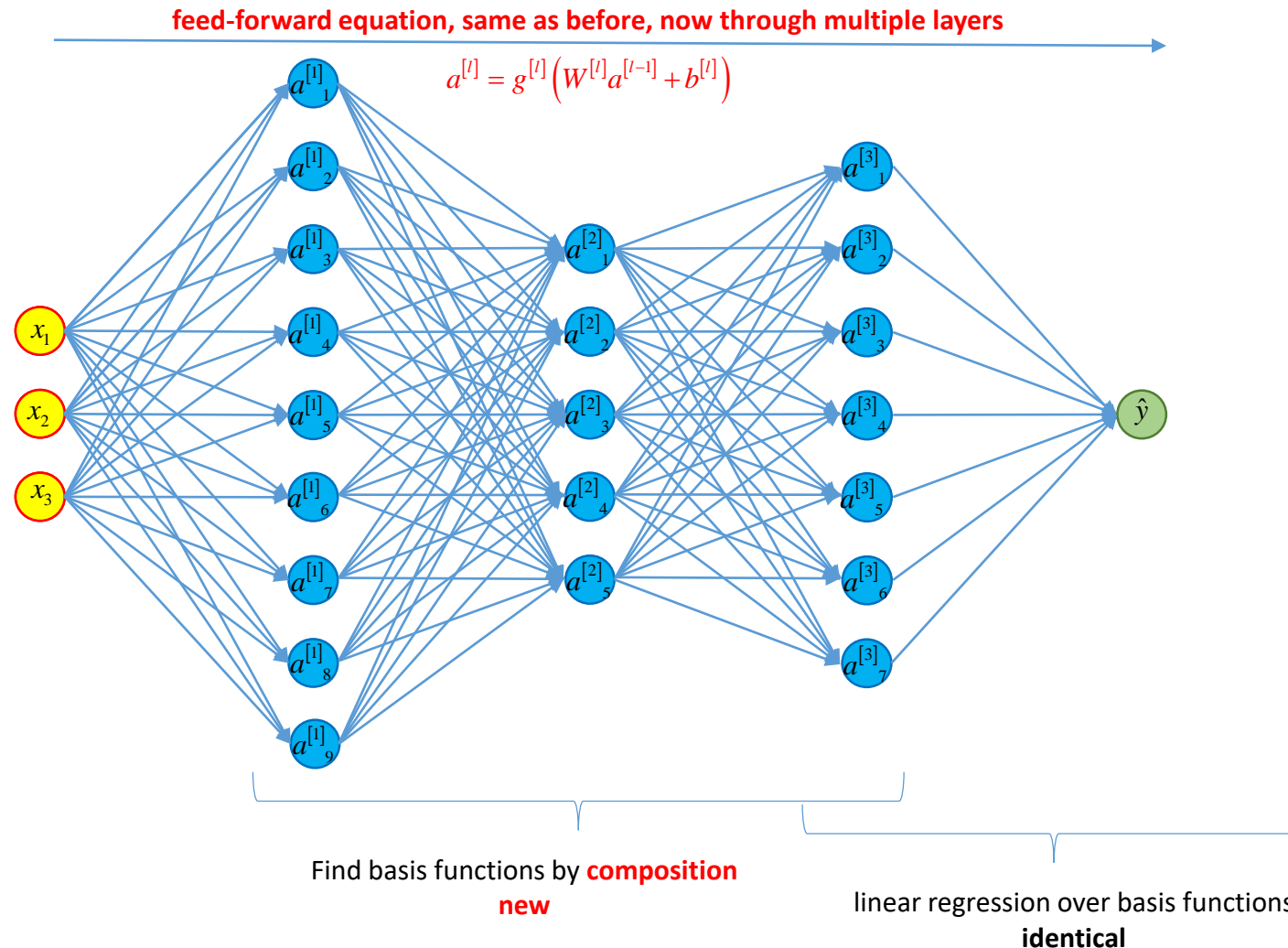
Universal Representation Theorem

- Polynomial regression offers mathematical guarantees
 - Combinations of monomials can approximate all smooth functions to arbitrary accuracy (with max degree p high enough)
 - Hence, polynomial regression (with enough monomials) can capture any (smooth) relation in the data
 - (At the cost of a high number of basis functions and the requirement of at least ten times that much data)
- Similarly, feed-forward ANNs with “perceptron” basis functions also offer mathematical guarantees
 - Known as the Universal Representation Theorem
 - States that a feed-forward perceptron may approximate any smooth function to arbitrary accuracy (with enough hidden units)
 - (Provided technical assumptions on activation functions, satisfied by all variants we presented)
 - See sketch of proof (with ReLU activation in dimension 1):
<https://www.quora.com/In-Machine-Learning-ReLU-activation-function-is-said-to-be-a-universal-function-approximation-can-you-show-just-how-it-is-possible-to-approximate-any-continuous-function-by-purely-combining-many-instances-of-ReLU/answer/Antoine-Savine-1>
- In conclusion
 - ANNs are superior to basis function regression in what they learn basis functions from the data
 - They scale to high dimension by regressing on a limited number of “most relevant” basis functions, found automatically
 - MLPs offer the same guarantees as basis function regression, also computationally friendly (matrix operations fit for parallel proc)
 - But they lose analytic solution for training and require iterative cost minimization
 - Iterative minimization procedures require gradients (all differentials) of the cost function on every iteration
 - Therefore a quick computation of derivatives is key to training in reasonable time

More on regression, basis functions and ANNs



Deep learning: multiple hidden layers



Basis function composition

- Shallow (one hidden layer) ANN:

- Prediction is a linear combination of hidden layer activations:
- Hidden layer activations are basis functions of the input layer:
- So we linearly regress on (learnt) basis functions of inputs

$$\hat{y} = W^{[2]}a^{[1]} + b^{[2]}$$
$$a^{[1]} = g^{[1]}(W^{[1]}x + b^{[1]})$$

- ANN with two hidden layers:

- Prediction is a linear combination of the **second** hidden layer:
- Second layer activations are basis functions of the first hidden layer:
- First hidden layer activations are basis functions of the input layer:
- So we regress on **functions of basis functions**

$$\hat{y} = W^{[3]}a^{[2]} + b^{[3]}$$
$$a^{[2]} = g^{[2]}(W^{[2]}a^{[1]} + b^{[2]})$$
$$a^{[1]} = g^{[1]}(W^{[1]}x + b^{[1]})$$

- Deep ANN with L layers (L-1 hidden layers)

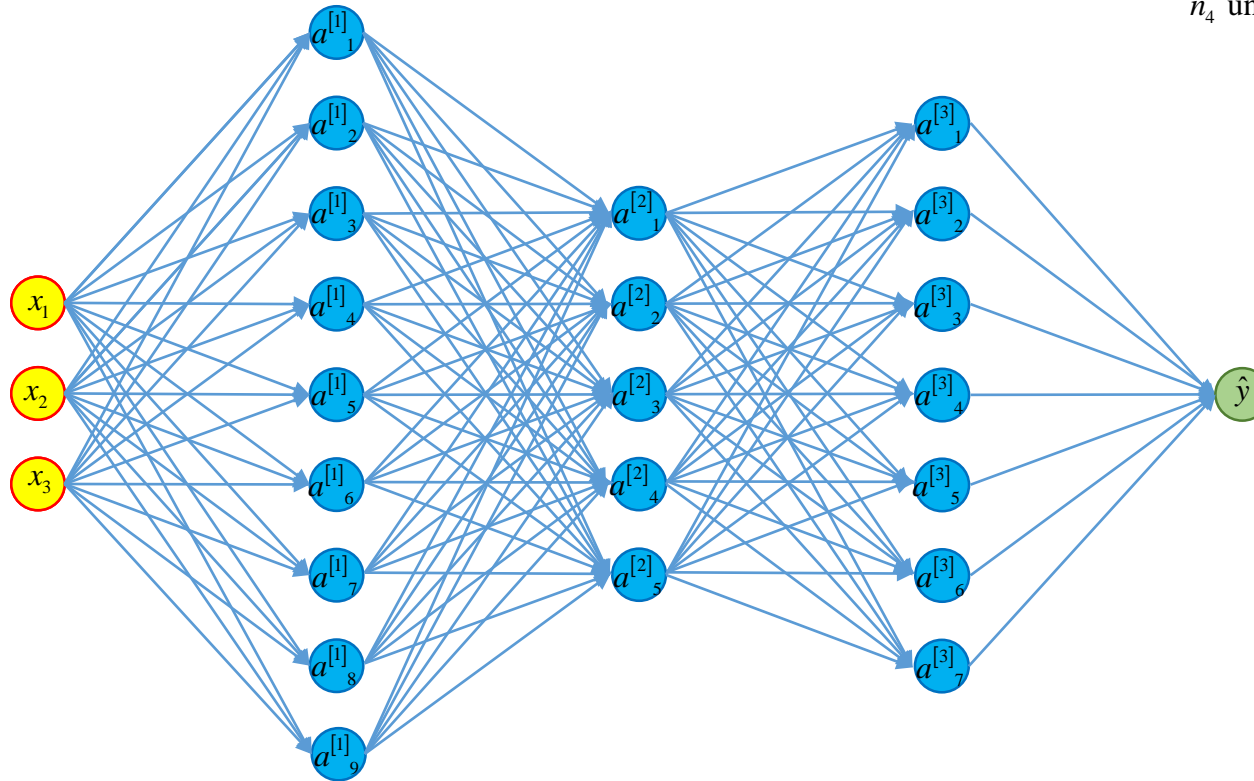
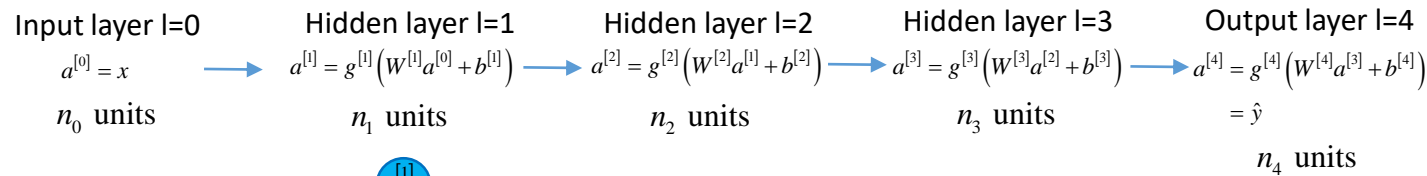
- Prediction is a linear combination of the **last** hidden layer:
- Last hidden layer contains basis functions of basis functions of basis functions ... of basis functions of the inputs from the feed-forward equation:
- So we regress on many **layers of composition** of basis functions

$$\hat{y} = W^{[L]}a^{[L-1]} + b^{[L]}$$
$$a^{[l]} = g^{[l]}(W^{[l]}a^{[l-1]} + b^{[l]})$$
$$a^{[0]} = x$$

Power of depth

- We know from the Universal Representation Theorem that:
 - An ANN with a single hidden layer can approximate any (smooth) function given enough units
 - So what is the point of multiple layers?
- In practice:
 - We have a finite -and limited- number of units
 - So we can only attain a limited sub-space of functions
- With composition, we attain a wider, richer space of functions
 - This intuition was demonstrated in limited cases and widely validated empirically
 - A convincing, general mathematical demonstration is (yet) unavailable
 - But a vast amount of empirical results confirm that deep ANNs much better represent target functions
 - In short, 5 layers of 10 units each attain a wider, richer sup-space of functions than 1 layer with 50 units
- This is the basis for deep learning
- The mechanics of deep ANNs is otherwise identical to shallow ANNs

Multi-layer perceptron (MLP)



feed-forward equation

$$a^{[l]} = g^{[l]}(W^{[l]}a^{[l-1]} + b^{[l]})$$

- Prediction: feed-forward equations

$$a^{[0]} = x, a^{[l]} = g^{[l]}(W^{[l]}a^{[l-1]} + b^{[l]}), y = a^{[L]}$$

- Parameters: for $1 \leq l \leq L$: $W^{[l]}(n_l \times n_{l-1}), b^{[l]}(n_l)$

- Number of parameters: $D = \sum_{l=1}^L n_l(1 + n_{l-1})$

- Complexity $\sim D = \sum_{l=1}^L n_l n_{l-1} = Ln^2$

- Training: find all $W^{[l]}$ and $b^{[l]}$ to minimise cost

$$C = \sum_{i=1}^m c_i, c_i = (\hat{y}^{(i)} - y^{(i)})^2$$

- Use iterative algorithms
compute all D differentials on each iteration

- Key: quickly compute
large number D of differentials

Summary and notations: general deep ANN

- Prediction

- Start with layer 0 = inputs $\begin{bmatrix} a^{[0]}_1 \\ \vdots \\ a^{[0]}_{n_0} \end{bmatrix} = \begin{bmatrix} x_1 \\ \vdots \\ x_{n_0} \end{bmatrix}$
- Feed-forward from layer l-1 to layer l $\rightarrow \begin{bmatrix} a^{[l]}_1 \\ \vdots \\ a^{[l]}_{n_l} \end{bmatrix} = \phi^{[l]} \left(\begin{bmatrix} a^{[l-1]}_1 \\ \vdots \\ a^{[l-1]}_{n_{l-1}} \end{bmatrix}; \begin{bmatrix} g^{[l]}_1 \\ \vdots \\ g^{[l]}_{D_l} \end{bmatrix} \right)$
- Prediction = layer L $\hat{y} = a^{[L]}_1$

- Parameters: $\begin{bmatrix} g^{[1]}_1 \\ \vdots \\ g^{[1]}_{D_1} \end{bmatrix}, \begin{bmatrix} g^{[2]}_1 \\ \vdots \\ g^{[2]}_{D_2} \end{bmatrix}, \dots, \begin{bmatrix} g^{[L-1]}_1 \\ \vdots \\ g^{[L-1]}_{D_{L-1}} \end{bmatrix}, \begin{bmatrix} g^{[L]}_1 \\ \vdots \\ g^{[L]}_{D_L} \end{bmatrix}$ so, ultimately: $\hat{y} = f \left(\begin{bmatrix} x_1 \\ \vdots \\ x_{n_0} \end{bmatrix}; \begin{bmatrix} g^{[1]}_1 \\ \vdots \\ g^{[1]}_{D_1} \end{bmatrix}, \begin{bmatrix} g^{[2]}_1 \\ \vdots \\ g^{[2]}_{D_2} \end{bmatrix}, \dots, \begin{bmatrix} g^{[L-1]}_1 \\ \vdots \\ g^{[L-1]}_{D_{L-1}} \end{bmatrix}, \begin{bmatrix} g^{[L]}_1 \\ \vdots \\ g^{[L]}_{D_L} \end{bmatrix} \right)$

- Training: find all the parameters theta

- To minimize sum of errors on training set with known labels $\min_{\begin{bmatrix} g^{[1]}_1 \\ \vdots \\ g^{[1]}_{D_1} \end{bmatrix}, \begin{bmatrix} g^{[2]}_1 \\ \vdots \\ g^{[2]}_{D_2} \end{bmatrix}, \dots, \begin{bmatrix} g^{[L-1]}_1 \\ \vdots \\ g^{[L-1]}_{D_{L-1}} \end{bmatrix}, \begin{bmatrix} g^{[L]}_1 \\ \vdots \\ g^{[L]}_{D_L} \end{bmatrix}} C = \sum_{i=1}^m c_i = (\hat{y}_i - y_i)^2$
- Takes iterative algorithm where all the differentials must be computed on every iteration: $\frac{\partial c_i}{\partial g^{[l]}_j}$ for all i, j and l

Summary and notations: MLP

- Special case of ANN with mathematical and computational benefits – by far the most common form of ANN

- Start with layer 0 = inputs $\begin{bmatrix} a_1^{[0]} \\ \vdots \\ a_{n_0}^{[0]} \end{bmatrix} = \begin{bmatrix} x_1 \\ \vdots \\ x_{n_0} \end{bmatrix}$

- Feed-forward from layer l-1 to layer l** $\rightarrow \begin{bmatrix} a_1^{[l]} \\ \vdots \\ a_{n_l}^{[l]} \end{bmatrix} = g^{[l]} \left(\begin{bmatrix} W_{1,1}^{[l]} & \dots & W_{1,n_{l-1}}^{[l]} \\ \vdots & \ddots & \vdots \\ W_{n_l,1}^{[l]} & \dots & W_{n_l,n_{l-1}}^{[l]} \end{bmatrix} \begin{bmatrix} a_1^{[l-1]} \\ \vdots \\ a_{n_{l-1}}^{[l-1]} \end{bmatrix} + \begin{bmatrix} b_1^{[l]} \\ \vdots \\ b_{n_l}^{[l]} \end{bmatrix} \right)$ where g is scalar, applied element-wise

- Prediction = layer L $\hat{y} = a_1^{[L]}$

- Parameters: all the $\begin{bmatrix} W_{1,1}^{[l]} & \dots & W_{1,n_{l-1}}^{[l]} \\ \vdots & \ddots & \vdots \\ W_{n_l,1}^{[l]} & \dots & W_{n_l,n_{l-1}}^{[l]} \end{bmatrix}$ and all the $\begin{bmatrix} b_1^{[l]} \\ \vdots \\ b_{n_l}^{[l]} \end{bmatrix}$ for all layers l

- Training: find all the parameters W and b

- Minimize sum of errors on training set with known labels $\min_{\begin{bmatrix} W_{1,1}^{[l]} & \dots & W_{1,n_{l-1}}^{[l]} \\ \vdots & \ddots & \vdots \\ W_{n_l,1}^{[l]} & \dots & W_{n_l,n_{l-1}}^{[l]} \end{bmatrix}, \begin{bmatrix} b_1^{[l]} \\ \vdots \\ b_{n_l}^{[l]} \end{bmatrix}, 1 \leq l \leq L} \left[C = \sum_{i=1}^m \{ c_i = (\hat{y}_i - y_i)^2 \} \right]$

- Takes iterative algorithm where all the differentials must be computed on every iteration: $\frac{\partial c_i}{\partial W_{j,k}^{[l]}}, \frac{\partial c_i}{\partial b_j^{[l]}}$ for all i, j, k and l

Back-Propagation

Elements of multivariate calculus

- We will derive the differentials of cost functions in general ANNs and MLPs
- Along with an algorithm to compute them extremely fast
- Since ANNs and MLPs are sequences of matrix operations: addition, product, element-wise activation, etc.
- We apply matrix differential calculus
- We briefly introduce a few elementary results, for a complete reference, see:
 - Mathematics for machine learning by Imperial College on Coursera
 - M. Giles, Extended collection of matrix derivative results (2008)
 - R. J. Barnes, Matrix Differentiation
 - B. Hage, AD of matrix calculations (public on ResearchGate)

Gradients

- Consider a scalar function of a vector: $f: \mathbb{R}^n \rightarrow \mathbb{R}$
- We call *gradient of f on x* and denote $\nabla f(x)$ the column vector of its derivatives: $\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(x) \\ \frac{\partial f}{\partial x_2}(x) \\ \dots \\ \frac{\partial f}{\partial x_n}(x) \end{bmatrix}$
- Key property of gradient:
the gradient vector gives the direction where the slope of f is steepest

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(x) \\ \frac{\partial f}{\partial x_2}(x) \\ \dots \\ \frac{\partial f}{\partial x_n}(x) \end{bmatrix}$$

derivatives in rows

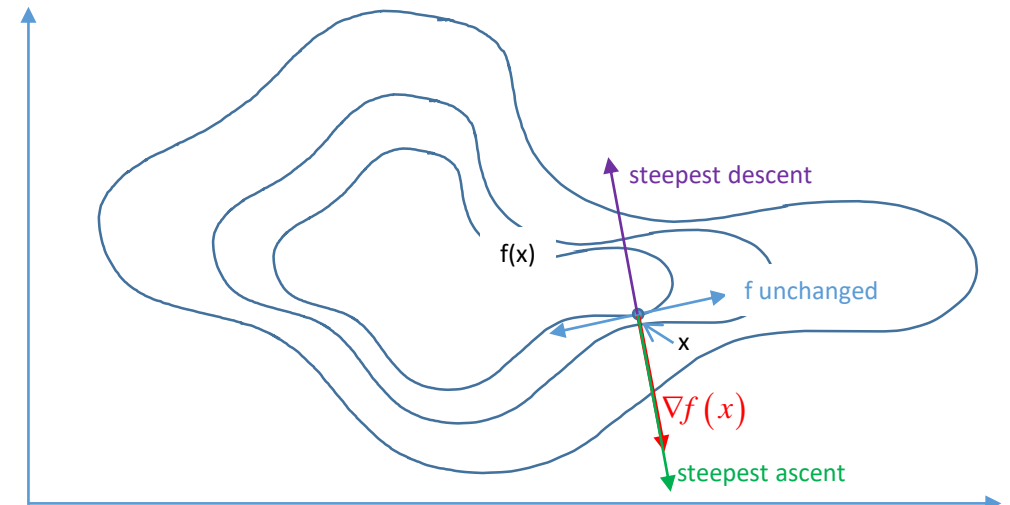
function in column

- Demonstration: to the 1st order

$$f(x + \Delta x) = f(x) + \sum_{i=1}^n \frac{\partial f}{\partial x_i} \Delta x_i = \nabla f(x)^T \Delta x = \langle \nabla f(x), \Delta x \rangle = \|\nabla f(x)\| \|\Delta x\| \cos(\nabla f(x), \Delta x)$$

Hence, for a fixed step size $\|\Delta x\|$

- steepest ascent is in the direction of $\nabla f(x)$ $\cos=1$
 - steepest descent is in the direction of $-\nabla f(x)$ $\cos=-1$
 - f is unchanged moving orthogonally to $\nabla f(x)$ $\cos=0$
- This key property is the basis of optimization algorithms and the process of training in machine learning

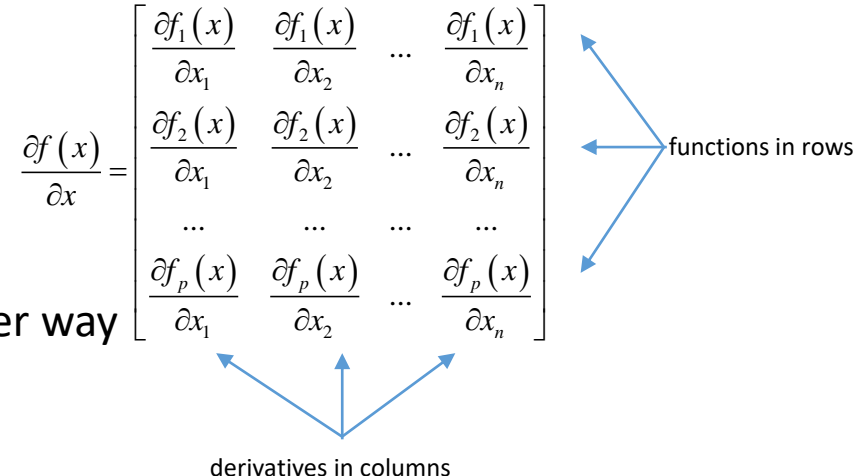


Matrix gradients

- The gradient of $f(x)$ has the same shape $(nx1)$ as x
- We want the same property scalar functions of a matrices: $f : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}$
- We denote $\frac{\partial f(M)}{\partial M}$ the matrix of its derivatives, of the same shape $(n \times p)$ as M

$$\frac{\partial f(M)}{\partial M} = \begin{bmatrix} \frac{\partial f(M)}{\partial M_{11}} & \frac{\partial f(M)}{\partial M_{12}} & \dots & \frac{\partial f(M)}{\partial M_{1p}} \\ \frac{\partial f(M)}{\partial M_{21}} & \frac{\partial f(M)}{\partial M_{22}} & \dots & \frac{\partial f(M)}{\partial M_{2p}} \\ \dots & \dots & \dots & \dots \\ \frac{\partial f(M)}{\partial M_{n1}} & \frac{\partial f(M)}{\partial M_{n2}} & \dots & \frac{\partial f(M)}{\partial M_{np}} \end{bmatrix}.$$

Jacobians

- Consider a vector function of a vector (p functions of x , x being of dim n) $f: \mathbb{R}^n \rightarrow \mathbb{R}^p, f(x) = \begin{bmatrix} f_1(x) \\ f_2(x) \\ \dots \\ f_p(x) \end{bmatrix}$
- The matrix $\frac{\partial f(x)}{\partial x}$ of the derivatives of the p functions to the n inputs is called Jacobian matrix
- Conventionally, Jacobians are often written with functions in rows and derivatives in columns: $\frac{\partial f(x)}{\partial x} = \begin{bmatrix} \frac{\partial f_1(x)}{\partial x_1} & \frac{\partial f_1(x)}{\partial x_2} & \dots & \frac{\partial f_1(x)}{\partial x_n} \\ \frac{\partial f_2(x)}{\partial x_1} & \frac{\partial f_2(x)}{\partial x_2} & \dots & \frac{\partial f_2(x)}{\partial x_n} \\ \dots & \dots & \dots & \dots \\ \frac{\partial f_p(x)}{\partial x_1} & \frac{\partial f_p(x)}{\partial x_2} & \dots & \frac{\partial f_p(x)}{\partial x_n} \end{bmatrix}$


functions in rows

derivatives in columns
- This is inconsistent with gradients, written the other way
Hence, confusing
- In this case, $\frac{\partial f(x)}{\partial x}$ is of shape $(p \times n)$
- And you can check that the chain rule is written: $\Rightarrow \underbrace{\frac{\partial h(x)}{\partial x}}_{q \times n} = \underbrace{\frac{\partial f(x)}{\partial g(x)}}_{q \times p} \underbrace{\frac{\partial g(x)}{\partial x}}_{p \times n}$ (outer function first)

$$g: \mathbb{R}^n \rightarrow \mathbb{R}^p$$

$$f: \mathbb{R}^p \rightarrow \mathbb{R}^q$$

$$h(x) = f[g(x)]: \mathbb{R}^n \rightarrow \mathbb{R}^q$$

Jacobians the other way

- We will write Jacobians the other way, consistently with gradients:
- This is consistent with gradients
 - Gradient = Jacobian of one-dimensional functions
 - Jacobian = horizontal concatenation of gradients

$$\frac{\partial f(x)}{\partial x} = \begin{bmatrix} \frac{\partial f_1(x)}{\partial x_1} & \frac{\partial f_2(x)}{\partial x_1} & \dots & \frac{\partial f_p(x)}{\partial x_1} \\ \frac{\partial f_1(x)}{\partial x_2} & \frac{\partial f_2(x)}{\partial x_2} & \dots & \frac{\partial f_p(x)}{\partial x_2} \\ \dots & \dots & \dots & \dots \\ \frac{\partial f_1(x)}{\partial x_n} & \frac{\partial f_2(x)}{\partial x_n} & \dots & \frac{\partial f_p(x)}{\partial x_n} \end{bmatrix}$$

derivatives in rows

functions in columns

- In this case, $\frac{\partial f(x)}{\partial x}$ is of shape $(n \times p)$
- And you can check that the chain rule is written:

$$\begin{aligned} g &: \mathbb{R}^n \rightarrow \mathbb{R}^p \\ f &: \mathbb{R}^p \rightarrow \mathbb{R}^q \\ h(x) &= f[g(x)]: \mathbb{R}^n \rightarrow \mathbb{R}^q \quad (\text{inner function first}) \\ \Rightarrow \underbrace{\frac{\partial h(x)}{\partial x}}_{n \times q} &= \underbrace{\frac{\partial g(x)}{\partial x}}_{n \times p} \underbrace{\frac{\partial f(x)}{\partial g(x)}}_{p \times q} \end{aligned}$$

Linear Jacobians

- If $f: \mathbb{R}^n \rightarrow \mathbb{R}^p$ is linear, i.e. $f(x) = Ax$ where A is a matrix of shape $(p \times n)$
- Then, you can convince yourself that its Jacobian matrix (with derivatives in rows) is the constant:

$$\frac{\partial f(x)}{\partial x} = A^T$$

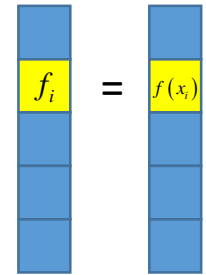
and not A which is maybe why people like Jacobians the other way around...

The diagram illustrates the equation $f_i = A_i * x$ using colored blocks. On the left, a vertical blue vector of 5 cells has its second cell highlighted in yellow and labeled f_i . This is followed by an equals sign. To the right of the equals sign is a 3x4 grid of blue cells, with its second row highlighted in yellow and labeled A_i . This is followed by an asterisk, and finally a vertical orange vector of 3 cells labeled x .

Element-wise Jacobians

- Consider a scalar function: $f: \mathbb{R} \rightarrow \mathbb{R}$

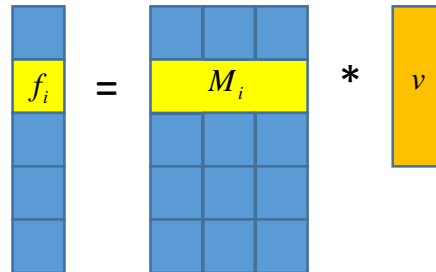
- We apply f element-wise to a vector x of dim n to obtain another vector of dim n : $f(x) = \begin{bmatrix} f(x_1) \\ f(x_2) \\ \dots \\ f(x_n) \end{bmatrix}$



- Then, the Jacobian of f is the $n \times n$ matrix: $\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} = f'(x_1) & \frac{\partial f_2}{\partial x_1} = 0 & \dots 0 \dots & 0 \\ \frac{\partial f_1}{\partial x_2} = 0 & \frac{\partial f_2}{\partial x_2} = f'(x_2) & \dots 0 \dots & 0 \\ 0 & 0 & \dots \frac{\partial f_i}{\partial x_i} = f'(x_i) \dots & \\ 0 & 0 & \dots 0 \dots & \frac{\partial f_n}{\partial x_n} = f'(x_n) \end{bmatrix} = \text{diag}[f'(x)]$

Vector functions of matrices

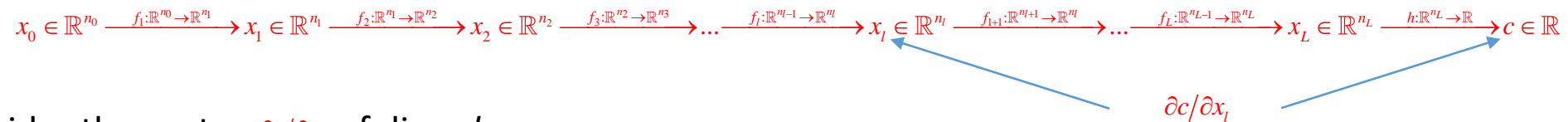
- Consider now a vector function of a matrix: $f: \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^q$
- In principle, its “derivative” is a 3D tensor of shape $(n \times p \times q)$: $\frac{\partial f}{\partial M} = \left(\frac{\partial f_k}{\partial M_{ij}} \right)$
- In the linear case where $f(M) = Mv$, v is a constant vector of size $p=q$
- We have $\frac{\partial f_i}{\partial M_i} = v^T$ and $\frac{\partial f_{j \neq i}}{\partial M_i} = 0$


$$\begin{bmatrix} f_i \end{bmatrix} = \begin{bmatrix} M_i \end{bmatrix} * \begin{bmatrix} v \end{bmatrix}$$

First touch with adjoints

- Consider the series of calculations:

- Initial input: x_0 vector of dim n_0
- Apply $x_1 = f_1(x_0)$ with $f_1: \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_1}$ (for avoidance of doubt, not necessarily element-wise) to compute a vector of dim n_1
- Apply $x_2 = f_2(x_1), x_3 = f_3(x_2), \dots, x_{l+1} = f_{l+1}(x_l), \dots, x_L = f_L(x_{L-1})$ to compute a sequence of vectors of dims $n_2, n_3, \text{ etc.}$
- Finally, apply a final scalar function to the final result of the chain: $c = h(x_L)$ with $h: \mathbb{R}^{n_L} \rightarrow \mathbb{R}$



- Consider the vector $\partial c / \partial x_i$ of dim n_i

- Means sensitivities of the final result c to small bumps in the components of x_i
 - All computations *before* stage i are fixed
 - All computations *after* stage i are impact by the bump to x_i

- Is called the adjoint of x_i and denoted \bar{x}_i

- Evidently, $\bar{x}_L = \nabla h(x_L)$ and by a straight application of the chain rule, we get the adjoint equation:

- Note that we have a formula for the last adjoint \bar{x}_L and for all adjoints \bar{x}_i function of the following adjoint \bar{x}_{i+1}

- Adjoint are therefore computed *in the reverse order*

$$\bar{x}_i = \frac{\partial c}{\partial x_i} = \frac{\partial x_{i+1}}{\partial x_i} \frac{\partial c}{\partial x_{i+1}} = \underbrace{\frac{\partial f_{i+1}(x_i)}{\partial x_i}}_{\substack{\text{Jacobian of } f_{i+1} \\ n_i \times n_{i+1}}} \bar{x}_{i+1}$$

adjoint of x_i adjoint of x_{i+1}

How to train your network, more formally

- **Recap:** we have an ANN of the general form implementing the following sequence of computations:

- Starting with inputs $a^{[0]} = x \in \mathbb{R}^{n_0}$
- Sequentially compute $a^{[l]} = \phi^{[l]}(a^{[l-1]}; \theta^{[l]}) \in \mathbb{R}^{n_l}$
- To make predictions $\hat{y} = a^{[L]} \in \mathbb{R}$

$$a^{[0]} = x \in \mathbb{R}^{n_0} \xrightarrow{\phi^{[1]}(\cdot; \theta^{[1]}): \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_1}} a^{[1]} \in \mathbb{R}^{n_1} \xrightarrow{\phi^{[2]}(\cdot; \theta^{[2]}): \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}} a^{[2]} \in \mathbb{R}^{n_2} \xrightarrow{\phi^{[3]}(\cdot; \theta^{[3]}): \mathbb{R}^{n_2} \rightarrow \mathbb{R}^{n_3}} \dots \xrightarrow{\phi^{[l]}(\cdot; \theta^{[l]}): \mathbb{R}^{n_{l-1}} \rightarrow \mathbb{R}^{n_l}} a^{[l]} \in \mathbb{R}^{n_l} \dots \xrightarrow{\phi^{[L]}(\cdot; \theta^{[L]}): \mathbb{R}^{n_{L-1}} \rightarrow \mathbb{R}^{n_L}} \hat{y} = a^{[L]} \in \mathbb{R}$$

- We also have a training set of m independent samples $x^{(m)}$ for the inputs x (each in dim n_0), along with corresponding labels $y^{(m)}$
- We want to train the ANN so its predictions correspond to the “best” predictions for y knowing x , that is, produce predictions $\hat{y} = E[y|x]$
- By definition, a conditional expectation is the closest function in the sense of least squares: $E[y|x] = \arg \min_f E[(y - f(x))^2]$
- And we know/have seen that:
 - (given enough height and depth) the ANN can represent any function f to arbitrary accuracy by adjusting its weights θ
 - (given enough training examples) $E[(y - f(x))^2] \approx \frac{1}{m} \sum_{i=1}^m [y^{(i)} - f(x^{(i)})]^2 = \frac{1}{m} \sum_{i=1}^m c_i$
- Therefore, “all we need to do” is find the weights θ^* that minimize the mean square prediction error on the training set: $\theta^* = \arg \min_{\theta} C(\theta) = \frac{1}{m} \sum_{i=1}^m c_i(\theta)$

Bad news good news

- There is no analytical solution for the optimal weights (unlike with linear regression)
- The cost $C(\theta)$ is a **non-convex** function of the parameters θ
- An algorithm guaranteed to find the minimum of a non-convex function **does not exist**, and probably never will
- So, the problem of training ANNs cannot, in theory, be solved
- This realization caused mathematicians to step away from deep learning in the 1990s
- Deep learning was revived recently by engineers along the lines of “this is not supposed to work, but let’s try anyway and see how we can take it”
- (And did it work!)
- The starting point is gradient descent:
 - Choose initial weights randomly
 - Take repeated steps
 - by a small amount called learning rate
 - along the direction of steepest descent, which we know is the direction opposite to the gradient
- Many heuristic improvements were made to improve and accelerate learning, but none is guaranteed to converge
- However, we understand the intuitions and take comfort in that they work remarkably well in a multitude of contexts

Differentials by finite differences

- Optimization algorithms
 - All learning algorithms start with a random guess and make steps according to the current gradient
 - We must compute and provide the gradient $\partial C / \partial \theta$ of the cost C to weights θ , repeatedly
 - Our goal is to compute it quickly
- Gradients by finite differences (FD):
 - First compute the cost C_0 , making predictions for all training examples
 - To compute the sensitivity of the cost to each (scalar) parameter θ_i :
 - Bump the parameter by a small amount ϵ
 - Re-compute the cost C_1 , repeating calculations for all training examples
 - $\partial C / \partial \theta_i \approx (C_1 - C_0) / \epsilon$
- FD is (obviously) linear in the (potentially immense) number of weights, therefore not viable
- But it has very desirable properties

Automatic derivatives with FD

- Implementation is straightforward with little scope for error
- Importantly, FD differentiation is **automatic**
 - Developers only write prediction (feed-forward) code
 - FD computes differentials automatically by calling the feed-forward code repeatedly
 - Developers don't need to write any differentiation code
- More importantly, FD automatically **synchronises** with modifications to feed-forward code
 - Tricks of the trade (see e.g. Coursera's Deep Learning Spec) to train deep nets faster and better:
 - Regularization: add a regularization term to cost function
 - Dropout: randomly drop connections between units
 - Batch Norm: normalize the mean and variance on all layers
 - And many more
 - All these modify feed-forward equations and code, and affect differentials
 - With manual differentiation, developers must consistently ensure consistency of evaluation and differentiation code
 - This is painful and prone to error
 - With automatic differentiation, differentiation always remains consistent with evaluation
- This being said, FD and other linear algorithms are not viable for training deep nets

Manual derivatives with AD

- Recall the sequence of calculations in the network:

$$a^{[0]} = x \in \mathbb{R}^{n_0} \xrightarrow{\phi^{[1]}(\cdot; \mathcal{G}^{[1]}): \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_1}} a^{[1]} \in \mathbb{R}^{n_1} \xrightarrow{\phi^{[2]}(\cdot; \mathcal{G}^{[2]}): \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}} a^{[2]} \in \mathbb{R}^{n_2} \xrightarrow{\phi^{[3]}(\cdot; \mathcal{G}^{[3]}): \mathbb{R}^{n_2} \rightarrow \mathbb{R}^{n_3}} \dots \xrightarrow{\phi^{[l]}(\cdot; \mathcal{G}^{[l]}): \mathbb{R}^{n_{l-1}} \rightarrow \mathbb{R}^{n_l}} a^{[l]} \in \mathbb{R}^{n_l} \dots \xrightarrow{\phi^{[L]}(\cdot; \mathcal{G}^{[L]}): \mathbb{R}^{n_{L-1}} \rightarrow \mathbb{R}^{n_L}} \hat{y} = a^{[L]} \in \mathbb{R} \rightarrow c = (\hat{y} - y)^2 \in \mathbb{R}$$

- We want to (quickly) compute all the $\frac{\partial C}{\partial \mathcal{G}^{[l]}}$

- Where the cost is the mean of the losses in the training set: $C = \frac{1}{m} \sum_{i=1}^m c_i$ so $\frac{\partial C}{\partial \mathcal{G}^{[l]}} = \frac{1}{m} \sum_{i=1}^m \frac{\partial c_i}{\partial \mathcal{G}^{[l]}}$

- We derive the derivatives of one loss c (dropping its index i in the training set) and sum them up in the end

- By the chain rule: $\frac{\partial c}{\partial \mathcal{G}^{[l]}} = \frac{\partial a^{[l]}}{\partial \mathcal{G}^{[l]}} \frac{\partial c}{\partial a^{[l]}} = \frac{\partial \phi^{[l]}(a^{[l-1]}; \mathcal{G}^{[l]})}{\partial \mathcal{G}^{[l]}} \frac{\partial c}{\partial a^{[l]}}$ and $\frac{\partial \phi^{[l]}(a^{[l-1]}; \mathcal{G}^{[l]})}{\partial \mathcal{G}^{[l]}}$ is **known**:

the shape of ϕ is a hyperparameter that defines the architecture of the network

the weights \mathcal{G} represent the current known state of the network

the activations a must be known, so a prior feed-forward evaluation is necessary before differentiation

- So we can first compute all the $\frac{\partial c}{\partial a^{[l]}}$, and then, apply the equation above to find all the $\frac{\partial c}{\partial \mathcal{G}^{[l]}}$

Manual derivatives with AD (2)

- Still from the chain rule, we see that: $\frac{\partial c}{\partial a^{[l]}} = \frac{\partial a^{[l+1]}}{\partial a^{[l]}} \frac{\partial c}{\partial a^{[l+1]}} = \frac{\partial \phi^{[l+1]}(a^{[l]}; \mathcal{G}^{[l+1]})}{\partial a^{[l]}} \frac{\partial c}{\partial a^{[l+1]}}$

$$a^{[0]} = x \in \mathbb{R}^{n_0} \xrightarrow{\phi^{[1]}(\cdot; \mathcal{G}^{[1]}): \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_1}} a^{[1]} \in \mathbb{R}^{n_1} \xrightarrow{\phi^{[2]}(\cdot; \mathcal{G}^{[2]}): \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}} a^{[2]} \in \mathbb{R}^{n_2} \xrightarrow{\phi^{[3]}(\cdot; \mathcal{G}^{[3]}): \mathbb{R}^{n_2} \rightarrow \mathbb{R}^{n_3}} \dots \xrightarrow{\phi^{[l]}(\cdot; \mathcal{G}^{[l]}): \mathbb{R}^{n_{l-1}} \rightarrow \mathbb{R}^{n_l}} a^{[l]} \in \mathbb{R}^{n_l} \dots \xrightarrow{\phi^{[L]}(\cdot; \mathcal{G}^{[L]}): \mathbb{R}^{n_{L-1}} \rightarrow \mathbb{R}^{n_L}} \hat{y} = a^{[L]} \in \mathbb{R} \rightarrow c = (\hat{y} - y)^2 \in \mathbb{R}$$

Note that this is the adjoint equation for the adjoints $\overline{a^{[l]}} = \frac{\partial c}{\partial a^{[l]}}$

- Therefore, we can calculate all the (adjoints) $\frac{\partial c}{\partial a^{[l]}}$:

$$\begin{aligned} \frac{\partial c}{\partial a^{[1]}} &= \frac{\partial a^{[2]}}{\partial a^{[1]}} \frac{\partial a^{[3]}}{\partial a^{[2]}} \frac{\partial a^{[4]}}{\partial a^{[3]}} \dots \frac{\partial a^{[L]}}{\partial a^{[L-1]}} \frac{\partial c}{\partial a^{[L]}} = \frac{\partial \phi^{[2]}(a^{[1]}; \mathcal{G}^{[2]})}{\partial a^{[1]}} \frac{\partial \phi^{[3]}(a^{[2]}; \mathcal{G}^{[3]})}{\partial a^{[2]}} \frac{\partial \phi^{[4]}(a^{[3]}; \mathcal{G}^{[4]})}{\partial a^{[3]}} \dots \frac{\partial \phi^{[L]}(a^{[L-1]}; \mathcal{G}^{[L]})}{\partial a^{[L-1]}} \frac{\partial c}{\partial a^{[L]}} \\ \frac{\partial c}{\partial a^{[2]}} &= \frac{\partial a^{[3]}}{\partial a^{[2]}} \frac{\partial a^{[4]}}{\partial a^{[3]}} \dots \frac{\partial a^{[L]}}{\partial a^{[L-1]}} \frac{\partial c}{\partial a^{[L]}} = \frac{\partial \phi^{[3]}(a^{[2]}; \mathcal{G}^{[3]})}{\partial a^{[2]}} \frac{\partial \phi^{[4]}(a^{[3]}; \mathcal{G}^{[4]})}{\partial a^{[3]}} \dots \frac{\partial \phi^{[L]}(a^{[L-1]}; \mathcal{G}^{[L]})}{\partial a^{[L-1]}} \frac{\partial c}{\partial a^{[L]}} \\ \frac{\partial c}{\partial a^{[3]}} &= \frac{\partial a^{[4]}}{\partial a^{[3]}} \dots \frac{\partial a^{[L]}}{\partial a^{[L-1]}} \frac{\partial c}{\partial a^{[L]}} = \frac{\partial \phi^{[4]}(a^{[3]}; \mathcal{G}^{[4]})}{\partial a^{[3]}} \dots \frac{\partial \phi^{[L]}(a^{[L-1]}; \mathcal{G}^{[L]})}{\partial a^{[L-1]}} \frac{\partial c}{\partial a^{[L]}} \\ &\dots \\ \frac{\partial c}{\partial a^{[l]}} &= \left[\prod_{i=l}^{L-1} \frac{\partial a^{[i+1]}}{\partial a^{[i]}} \right] \frac{\partial c}{\partial a^{[L]}} = \left[\prod_{i=l}^{L-1} \frac{\partial \phi^{[i+1]}(a^{[i]}; \mathcal{G}^{[i+1]})}{\partial a^{[i]}} \right] \frac{\partial c}{\partial a^{[L]}} \end{aligned}$$

where we recall that all the $\frac{\partial \phi^{[i+1]}(a^{[i]}; \mathcal{G}^{[i+1]})}{\partial a^{[i]}}$ are known, we “just” need to compute their matrix products

Manual derivatives with AD (3)

- We have our formula: $\frac{\partial c}{\partial a^{[l]}} = \left[\prod_{i=l}^{L-1} \frac{\partial \phi^{[i+1]}(a^{[i]}; \mathcal{G}^{[i+1]})}{\partial a^{[i]}} \right] \frac{\partial c}{\partial a^{[L]}} = \left[\prod_{i=l}^{L-1} \frac{\partial \phi^{[i+1]}(a^{[i]}; \mathcal{G}^{[i+1]})}{\partial a^{[i]}} \right] \frac{\partial c}{\partial \hat{y}}$
- And we know that $c = (\hat{y} - y)^2$ so $\frac{\partial c}{\partial a^{[L]}} = \frac{\partial c}{\partial \hat{y}} = 2(\hat{y} - y)$ “two prediction errors”, which we recall is a real number
- We must compute $\prod_{i=l}^{L-1} \frac{\partial \phi^{[i+1]}(a^{[i]}; \mathcal{G}^{[i+1]})}{\partial a^{[i]}}$ for every layer l , where we recall that all the $\frac{\partial \phi^{[i+1]}(a^{[i]}; \mathcal{G}^{[i+1]})}{\partial a^{[i]}}$ are known from:
 - The architecture of the network, which gives the form all the vector functions ϕ , hence, also their Jacobians
 - The parameters θ on the current iteration
 - The activations a computed while making predictions
(prior to computing derivatives, we must compute the predictions and store all the activations on all the layers)
- What is the most efficient way to compute this matrix product $\prod_{i=l}^{L-1} \frac{\partial \phi^{[i+1]}(a^{[i]}; \mathcal{G}^{[i+1]})}{\partial a^{[i]}}$ knowing all the matrices in the product?

Manual derivatives with AD (4)

- For example, let us compute the derivatives to the 1st layer: $\frac{\partial c}{\partial a^{[1]}} = \frac{\partial a^{[2]}}{\partial a^{[1]}} \frac{\partial a^{[3]}}{\partial a^{[2]}} \frac{\partial a^{[4]}}{\partial a^{[3]}} \cdots \frac{\partial a^{[L-1]}}{\partial a^{[L-2]}} \frac{\partial a^{[L]}}{\partial a^{[L-1]}} \frac{\partial c}{\partial a^{[L]}}$
- We could compute it left to right, layer 1 to layer L , in the order of the feed-forward equations
- Or, we could compute it right to left, layer L to layer 1, in the reverse order
- Note that $a^{[L]} = \hat{y}$ and $c = (\hat{y} - y)^2$ are real numbers, so the rightmost term $\frac{\partial c}{\partial a^{[L]}} = \frac{\partial c}{\partial \hat{y}} = 2(\hat{y} - y)$ is also a real number

and the second rightmost term $\frac{\partial a^{[L]}}{\partial a^{[L-1]}} = \frac{\partial \hat{y}}{\partial a^{[L-1]}}$ is a vector of dim n_{L-1} , the gradient of the prediction to the basis functions in the final hidden layer

- All the other terms $\frac{\partial a^{[l]}}{\partial a^{[l-1]}}$, $2 \leq l < L$ are matrices of shape $n_{l-1} \times n_l$
- Crucially, the result $\frac{\partial c}{\partial a^{[1]}}$ is a vector (of dim n_1) because the loss c is scalar

Manual derivatives with AD (5)

- So:

$$\underbrace{\begin{bmatrix} \frac{\partial c}{\partial a^{[1]}} \end{bmatrix}}_{\text{vector}} = \underbrace{\begin{bmatrix} \frac{\partial a^{[2]}}{\partial a^{[1]}} \end{bmatrix}}_{\text{matrix}} \underbrace{\begin{bmatrix} \frac{\partial a^{[3]}}{\partial a^{[2]}} \end{bmatrix}}_{\text{matrix}} \underbrace{\begin{bmatrix} \frac{\partial a^{[4]}}{\partial a^{[3]}} \end{bmatrix}}_{\text{matrix}} \cdots \underbrace{\begin{bmatrix} \frac{\partial a^{[L-1]}}{\partial a^{[L-2]}} \end{bmatrix}}_{\text{matrix}} \underbrace{\begin{bmatrix} \frac{\partial a^{[L]} \frac{\partial c}{\partial a^{[L]}} \end{bmatrix}}_{\text{vector}}$$

→ forward calculation
← backward calculation

- Forward calculation:
 - Start with a matrix
 - Repeatedly multiply matrices by matrices to obtain matrices, which is of cubic complexity
 - Collapse to a vector by right multiplication with a vector on the final hidden layer
- Backward calculation:
 - Start with a **vector**
 - Repeatedly multiply matrices by vectors to obtain vectors, which is of quadratic complexity
- Backward calculation is one order of magnitude faster!

Back-propagation algorithm

- Back-prop: algorithm to compute derivatives in the reverse order, one order of magnitude faster

- Algorithm:

1. Compute prediction by feed-forward through the network, storing all hidden activations a^l
2. Compute the adjoints $\overline{a^l} = \partial c / \partial a^l$ of all the layer activations in the reverse order

a. Start with $\overline{a^L} = \frac{\partial c}{\partial a^L} = \frac{\partial c}{\partial \hat{y}} = 2(\hat{y} - y)$

b. Recursively compute $\overline{a^l} = \frac{\partial c}{\partial a^l} = \frac{\partial a^{l+1}}{\partial a^l} \frac{\partial c}{\partial a^{l+1}} = \frac{\partial \phi^{l+1}(a^{l+1}; g^{l+1})}{\partial a^{l+1}} \frac{\partial c}{\partial a^{l+1}} = \underbrace{\frac{\partial \phi^{l+1}(a^{l+1}; g^{l+1})}{\partial a^{l+1}}}_{\text{known}} \overline{a^{l+1}}$ from layer $L-1$ to layer 1

Note we get all the $\overline{a^l} = \partial c / \partial a^l$ in a single backward sweep, computing them all from one another

Also recall the Jacobian matrices $\frac{\partial \phi^{l+1}(a^{l+1}; g^{l+1})}{\partial a^{l+1}}$ are known from the architecture, current weights and pre-calculated activations

3. Compute the desired derivatives $\frac{\partial c}{\partial g^l} = \frac{\partial a^l}{\partial g^l} \frac{\partial c}{\partial a^l} = \frac{\partial \phi^l(a^l; g^l)}{\partial g^l} \frac{\partial c}{\partial a^l} = \frac{\partial \phi^l(a^l; g^l)}{\partial g^l} \overline{a^l}$ where, again, $\frac{\partial \phi^l(a^l; g^l)}{\partial g^l}$ is known

Back-propagation = AD

- The core of back-prop is the recursive reverse equation: $\overline{a^{[l]}} = \frac{\partial a^{[l+1]}}{\partial a^{[l]}} \overline{a^{[l+1]}}$
- We have seen this equation before and called it adjoint equation
- Back-prop is the application of the adjoint equation to neural nets
- **Back-prop (also called AD) is not limited to neural nets**
- **It is applicable to any calculation as long as:**
 - The final result is scalar, so adjoints are vectors and not matrices
 - The computation is organized in layers, where every layer is a known function of the previous layer only
 - We will see that any computation can be organized in layers, this is called the computation graph
- So back-prop (AD) is applicable to any computation, **as long as the final result is scalar**
- And it always computes all adjoints (derivatives) in a single (backward) sweep through the calculation graph

$$a^{[0]} = x \in \mathbb{R}^{n_0} \xrightarrow{\phi^{[1]}(\cdot; g^{[1]}): \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_1}} a^{[1]} \in \mathbb{R}^{n_1} \xrightarrow{\phi^{[2]}(\cdot; g^{[2]}): \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}} a^{[2]} \in \mathbb{R}^{n_2} \xrightarrow{\phi^{[3]}(\cdot; g^{[3]}): \mathbb{R}^{n_2} \rightarrow \mathbb{R}^{n_3}} \dots \xrightarrow{\phi^{[l]}(\cdot; g^{[l]}): \mathbb{R}^{n_{l-1}} \rightarrow \mathbb{R}^{n_l}} a^{[l]} \in \mathbb{R}^{n_l} \dots \xrightarrow{\phi^{[L]}(\cdot; g^{[L]}): \mathbb{R}^{n_{L-1}} \rightarrow \mathbb{R}^{n_L}} \hat{y} = a^{[L]} \in \mathbb{R} \rightarrow c = (\hat{y} - y)^2 \in \mathbb{R}$$

Back-propagation through MLPs

- MLP is a particular case where:

$$\phi^{[l]}(a^{[l-1]}; W^{[l]}, b^{[l]}) = g^{[l]}(W^{[l]}a^{[l-1]} + b^{[l]})$$

scalar activation function
for layer l
applied element-wise

matrix of weights
for layer l

vector of biases
for layer l

- It is convenient to split each layer in two: a linear combination followed by an activation:

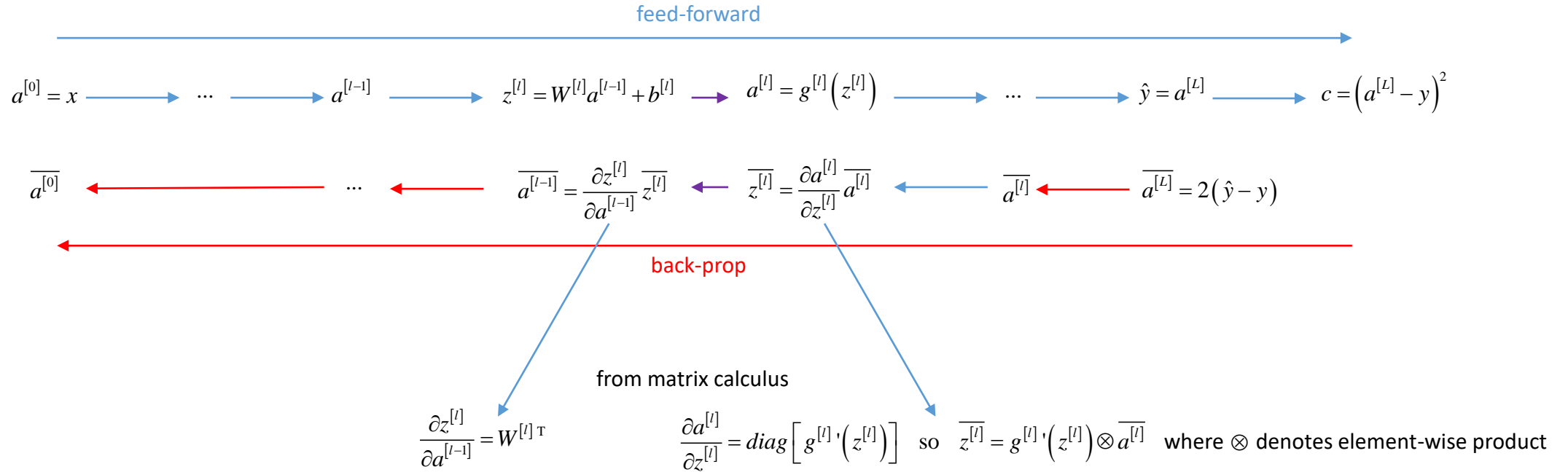
$$(1) \quad z^{[l]} = W^{[l]}a^{[l-1]} + b^{[l]} \quad \longrightarrow \quad (2) \quad a^{[l]} = g^{[l]}(z^{[l]})$$

- And we have the adjoint equations in the reverse order:

$$(1) \quad \overline{a^{[l-1]}} \equiv \frac{\partial c}{\partial a^{[l-1]}} = \frac{\partial z^{[l]}}{\partial a^{[l-1]}} \frac{\partial c}{\partial z^{[l]}} \equiv \frac{\partial z^{[l]}}{\partial a^{[l-1]}} \overline{z^{[l]}} \quad \longleftarrow \quad (2) \quad \overline{z^{[l]}} \equiv \frac{\partial c}{\partial z^{[l]}} = \frac{\partial a^{[l]}}{\partial z^{[l]}} \frac{\partial c}{\partial a^{[l]}} \equiv \frac{\partial a^{[l]}}{\partial z^{[l]}} \overline{a^{[l]}}$$

- Where $\frac{\partial z^{[l]}}{\partial a^{[l-1]}}$ and $\frac{\partial a^{[l]}}{\partial z^{[l]}}$ are given by a direct application of matrix calculus (see next slide)

Back-propagation through MLPs (2)



and finally (also from matrix calculus) we have the sensitivities to weights from the sensitivities to activation:

$$\frac{\partial c}{\partial W^{[l]}} = \bar{z}^{[l]} a^{[l-1] \top}$$

$$\frac{\partial c}{\partial b^{[l]}} = \frac{\partial z^{[l]}}{\partial b^{[l]}} \frac{\partial c}{\partial z^{[l]}} = \frac{\partial c}{\partial z^{[l]}} = \bar{z}^{[l]}$$

because:
$$\frac{\partial c}{\partial W_i^{[l]}} = \frac{\partial z_i^{[l]}}{\partial W_i^{[l]}} \frac{\partial c}{\partial z_i^{[l]}} = \frac{\partial z_i^{[l]}}{\partial W_i^{[l]}} \bar{z}_i^{[l]} = \bar{z}_i^{[l]} a^{[l-1] \top}$$

MLP back-propagation algorithm

- Which gives us the MLP back-prop, generally described in literature and implemented in frameworks
- Algorithm:
 1. Compute prediction by feed-forward through the network, storing all hidden combinations z^l and activations a^l
 2. Compute the adjoints $\overline{a^{[l]}} = \partial c / \partial a^{[l]}$ and $\overline{z^{[l]}} = \partial c / \partial z^{[l]}$ of all the layer combinations and activations in the reverse order
 - a. Start with $\overline{a^{[L]}} = 2(\hat{y} - y)$
 - b. Recursively compute $\overline{z^{[l]}} = g^{[l]'}(z^{[l]}) \otimes \overline{a^{[l]}}$ and then $\overline{a^{[l-1]}} = W^{[l] \top} \overline{z^{[l]}}$ from layer $L-1$ to layer 1
 3. Compute the desired derivatives $\frac{\partial c}{\partial W^{[l]}} = \overline{z^{[l]}} a^{[l-1] \top}$ and $\frac{\partial c}{\partial b^{[l]}} = \overline{z^{[l]}}$ of loss to weights and biases

Back-propagation: complexity

- Recall the leading complexity of feed-forward is a matrix by vector product by layer of quadratic complexity in the number of units
- We need a prior feed-forward before back-prop to compute and store all the z^l and a^l
- The back-prop equations are executed once on every layer (in reverse order) to compute derivatives to activations and weights:
 - $\overline{z^l} = g^{[l]'}(z^l) \otimes \overline{a^l}$ is of linear complexity (element-wise multiplication of two vectors)
 - $\overline{a^{l-1}} = W^{[l]T} \overline{z^l}$ is of quadratic complexity (matrix by vector product)
 - $\frac{\partial c}{\partial W^{[l]}} = \overline{z^l} \overline{a^{l-1}T}$ is of quadratic complexity
 - $\frac{\partial c}{\partial b^{[l]}} = \overline{z^l}$ is immediate

Complexity and memory considerations

- So the complexity of a complete back-prop to compute all derivatives is 3 times the complexity of a feed-forward prediction:
 - One necessary feed-forward prediction to compute and store all the z_l and a_l
 - The back-propagation itself takes one matrix by vector product, quadratic in units, like feed-forward (but in reverse order)
 - Finally, the computation of derivatives to weights adds another matrix by vector product on each layer
- As promised, we compute a large number of derivatives in constant time
 - The actual back-prop is as fast as one prediction (the other 2 are for the prior prediction and derivatives to weights)
 - Independently on the number of derivatives
- It works because:
 - The end result is scalar, so its derivatives to vectors are vectors
 - Reverse adjoint propagation collapses matrix by vector products into vectors, so we have quadratic complexity, like feed-forward
 - Derivatives are computed analytically in one sweep, there is no need to repeatedly bump many parameters one by one
- It consumes significant memory
 - All combinations and activations must be stored for all layers
 - In addition to all the weights and biases

Automatic backprop?

- Reverse adjoint propagation
 - Also called back-propagation or back-prop in ML lingo, or **adjoint differentiation** (AD) in general
 - Computes all differentials of the cost to the many parameters in a deep net in a time similar to one evaluation of the cost
 - Achieves this remarkable result by reversing the order of the calculations in the differentiation
 - Is the only viable means to train deep nets in reasonable time
 - Relies on the cost being scalar, hence its adjoints are vectors
 - **Applies to any scalar calculation, not just neural nets**
- However manual AD code is
 - Complicated and prone to error
 - Painful to maintain and synchronize as we modify feed-forward code
- AAD: **automatic** adjoint differentiation
 - Next, we learn to **automate** AD
 - Since AD propagates adjoints in the reverse order through evaluation graphs
 - We will learn to automatically generate graphs so we can back-prop derivatives through them
 - Not only for deep nets, but for any calculation
 - Including large Monte-Carlo simulations
- First, we illustrate the application of deep learning to financial derivatives

Deep learning in finance

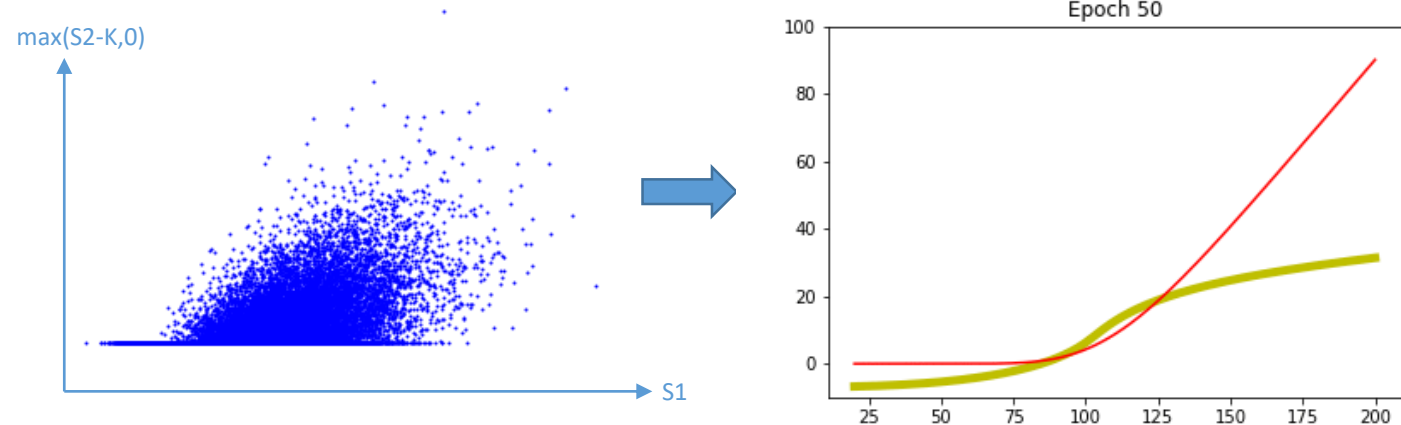
Example: learn Black-Scholes from data

- We reuse the same example as for linear and quadratic regression

- Simulate data for the spot in 1y and a call payoff in 2y

- Train a neural net to find the future price of the call in 1y as a function of the spot prevailing in 1y

- See if it can learn Black & Scholes from data alone



- The net learns Black & Scholes from the data cloud alone
- The red line (Black & Scholes's formula) is only for the viewer's reference
- The net doesn't know anything about Black & Scholes at any stage

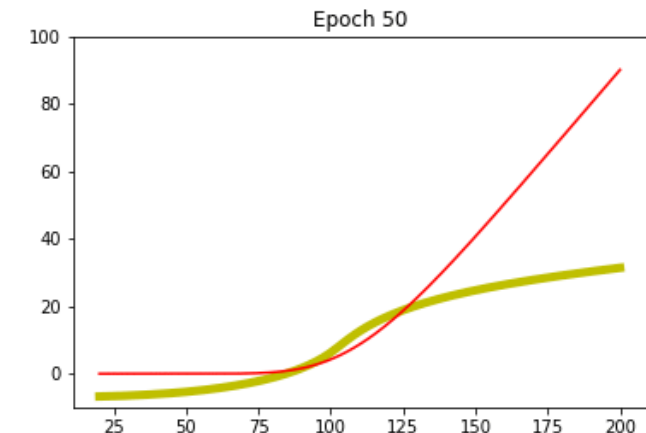
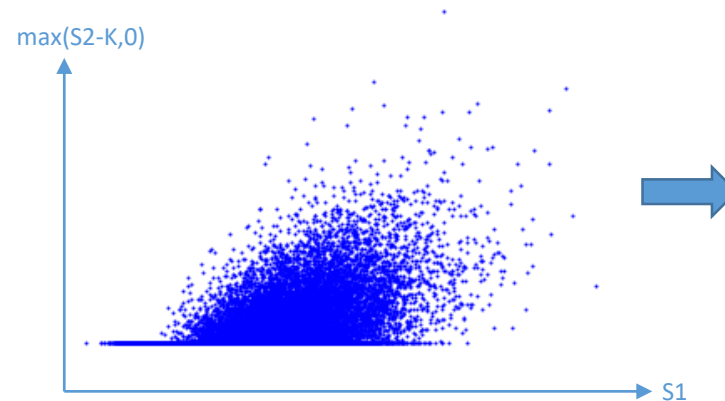
Simple ANNs with TensorFlow

- TensorFlow
 - Machine learning framework open sourced by Google
 - The most popular framework for machine learning and deep learning
- What is TensorFlow exactly?
 - Another math/numerical/matrix library with API in Python
 - With 3 perks (in increasing order of perkness)
 1. High level functions to work with ANNs easily and conveniently
 2. Matrix calculus executed on multi-core CPU and GPU
 3. Built-in back-prop so TF computes derivatives of your code, in constant time, behind the scenes (a form of AAD)
- Working with TensorFlow
 - To do all this, TF needs to know computation graphs: all computations with order and dependency, before they are executed
 - Hence, working with TF is surprising at first, you must:
 - First “register” the calculations with TF, that is, build the graph in TF’s memory space
 - Then, ask TF to run the graph, and back-prop through the graph when derivatives are needed

Simple example walkthrough

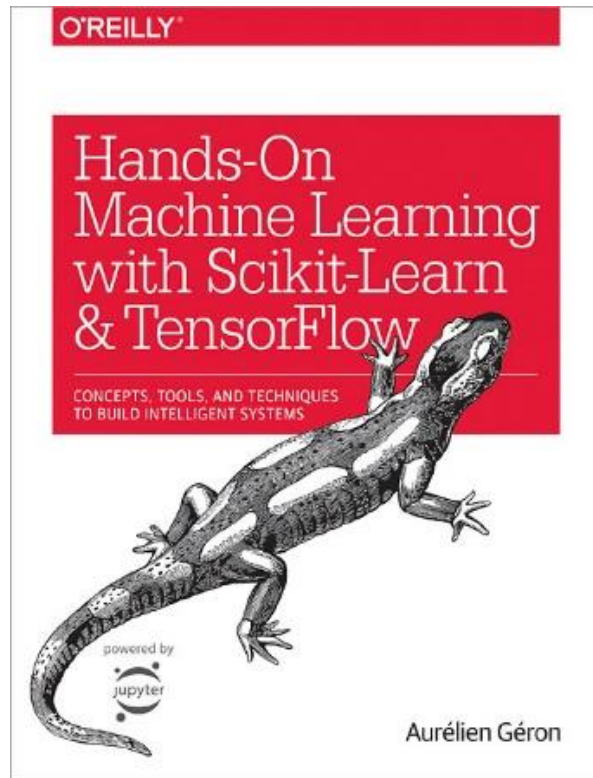
- Head to <http://GitHub/asavine/CompFinance/Workshop/>

- Run the notebook **dlBlackScholes.ipynb**
(TensorFlow must be installed)

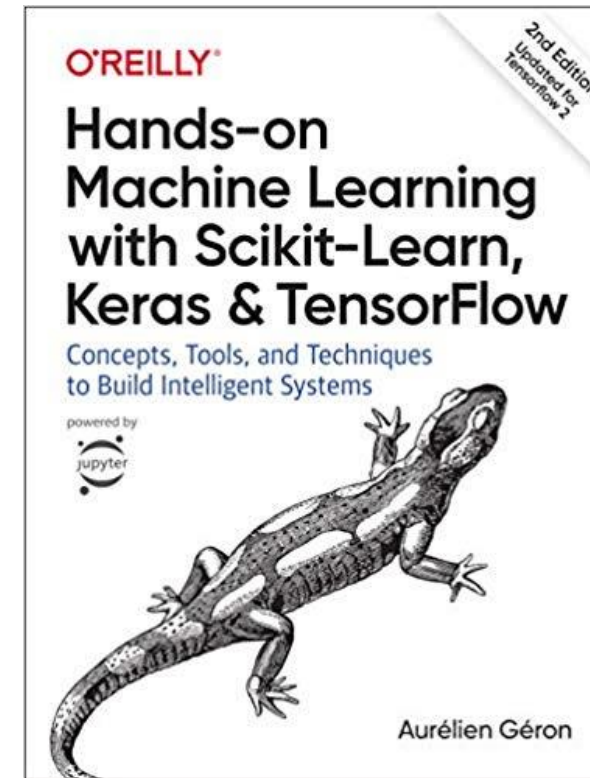


- This notebook demonstrates
 - The (basic) operation of TensorFlow
 - How to build and train (simple) neural nets in practice
 - Application to financial derivatives (revaluation)
 - How to implement a simple example in Black & Scholes

Learning TensorFlow



Summer 2019



Revaluation nets

- What we just saw is a (very simple) example of a revaluation problem:
What is the future value in simulated scenarios of a derivatives transaction or book of transactions?
- Vast domain of application:
 - Pricing callable options and exotics with Monte-Carlo
 - Simulate state variables up to the exercise date
 - Use the net to find the continuation value on the exercise date
 - Determine if the call is exercised or not in this scenario
 - Extension of the well-known “Least Square Method” of Longstaff-Schwartz and Carriere (late 1990s) to modern deep learning
 - More importantly, regulatory simulations like XVA, MVA, CCR or FRTB
 - Simulate the future value of derivatives books in a large number of scenarios to compute exposures
 - Normally require nested simulations (MC to estimate book value inside MC to simulate market to exposure date!)
 - Revaluation nets offer alternative orders of magnitude faster and more accurate
 - Read the intro of the notebook in <http://GitHub/asavine/CompFinance/Workshop/dlBlackScholes.ipynb>

Revaluation nets (2)

- Revaluation nets are disposable
 - Trained for a given set of transactions in a given model, calibrated to current market data
 - Market, model and transaction parameters get encoded in the weights of the network
 - Only input is simulated state vector
 - The input is only valid with the market, model and transaction it was trained with
 - Revaluation nets are trained once, used once
 - Hence, training must be very fast (seconds)
- Revaluation nets are very general, work with models and transactions of any complexity
- Revaluation nets resolve key problems with regulatory simulations, in particular XVA, FRTB, CCR and MVA
- Upcoming article “Deep Analytics” by Brian Huge and Antoine Savine

Pricing nets

- An apparently similar problem is pricing
- Context:
 - Given a dynamic model with a number of parameters
 - Current market: asset prices, interest rates, etc.
 - Dynamic parameters: volatilities, volatilities of volatilities, correlations, etc.
 - Not an exotic model a la Dupire with vast (term structures) of parameters
 - But a European model a la SABR/Heston/RoughBergomi with a limited number of parameters
 - And some transaction, also with a number of parameters
 - Generally, European option or first order exotic (barrier) with few parameters like expiry, strike, barrier
- Pricing problem:
 - Only simplest models like Black & Scholes have explicit closed form solutions
 - More realistic models require
 - Expansions/approximations (SABR/ZABR)
 - Numerical integration (Heston)
 - Multi-Factor FDM
 - Slow Monte-Carlo simulations (RoughVol)

Pricing nets (2)

- Can we do better with a neural net?
 - Generate a vast training set, sampling the space of all market, model and transaction parameters
 - For each sampled parameter set, produce a price, e.g. with Monte-Carlo simulations
 - Note the training set alone may take hours or days
 - Train a neural net to generalize pricing from the training set
 - Apply the trained net as a very fast, very accurate ~almost closed form~ pricer
- Train once, use forever
- Very different from reval nets
 - May train for hours or days, forever reusable
 - Weights only encode pricing mechanics. Market, dynamic and transaction parameters are inputs
 - Training set is sampled, not simulated - sampling it correctly is perhaps the major challenge
 - One net per model (e.g. McGhee's SABRNet for SABR or Horvath and al's RoughNet for RoughBergomi)
 - Limited to European models and transactions, not likely to extend to exotics

Pricing nets (3)

- Recent results
 - McGhee (2018) built a SABR net, more accurate than Hagan's approximation and orders of magnitude faster than 2D FDM
 - Ferguson and Greene (also 2018) built a net for pricing basket options with parameters (inputs) including correlation
 - Horvath and al. (2019) built a net for pricing with Rough Volatility
 - This is particularly significant
 - Because there exist no alternatives for pricing with rough vol other than slow Monte-Carlo
 - RoughNet therefore opens new applications for calibration, pricing and risk management with rough vol
 - So far, rough vol was restricted to academic circles due to absence of fast pricing, it may now enter production

Deeper learning for financial derivatives

	Reval nets	Pricing nets
Output	Future price depending on scenario	Current price
Inputs	Simulated state variables	Market, model and transaction parameters
Training set	State variables to payoffs = realizations per scenario	Model parameters to prices = expectations across scenarios, per parameter set
Weight encoding	Valuation mechanics + market, model and transaction parameters	Valuation mechanics only
Applications	Callable exotics with MC XVA, MVA, FRTB, CCR, etc.	Europeans and first order exotics Stochastic volatility/rough vol/jump/etc. Calibration/pricing/risk management
Reuse	Disposable	Reusable
Model and transaction scope	Any complexity	Restricted to Europeans and simple exotics
Accuracy requirement	Medium	High
Training time	Seconds	Hours
Reference nets	DeepAnalytics (Huge & Savine, coming 2019)	McGhee's SABRNet (2018) Ferguson & Greene's BasketNet (2018) Horvath and al's RoughNet (2019)

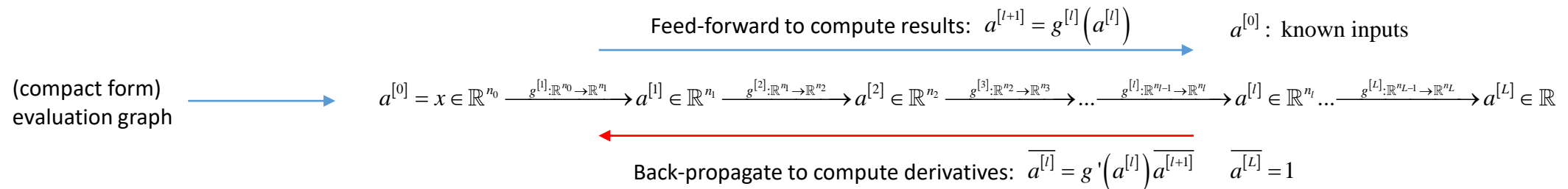
Deeper learning

- We studied in detail the fundamental machinery of feed-forward and back-propagation
- There is a lot more to deep learning:
 - Speed and accuracy improvements: regularization, dropout, batch-norm, stochastic gradient descent, etc.
 - Convolutional nets with many impressive applications in computer vision
 - Recurrent nets with many impressive applications in natural language processing (NLP)
And also time series forecasting, proprietary trading strategies, etc.
 - Reinforcement learning, leveraging DL to learn strategies for playing games, driving cars or trading currencies
- Deep learning is here to stay
 - Overhyped due to recent successes in computer vision, NLP and chess/go --- hype will fade
 - This being said, DL genuinely offers efficient, elegant solutions to otherwise intractable problems in finance and elsewhere
 - DL will not replace financial models but it will be part of the financial quant toolbox
 - Financial quants and derivatives professionals cannot ignore it

AD through evaluation graphs

Evaluation graphs and adjoint propagation

- Evaluation graph
 - We have seen evaluation graphs for neural nets
 - Evaluation graphs express the order and dependencies of “atomic” calculations involved in a larger evaluation
 - Any evaluation defines a graph
 - Hence, every evaluation defines a sequence of feed-forward operations, similarly to a neural net
 - Chapter 9 of the Modern Computational Finance book builds evaluation graphs in memory with operator overloading in C++
- Adjoint Differentiation is back-propagation over graphs
 - Not limited to neural nets, applicable to any evaluation
 - Propagates adjoints backwards through the evaluation graph
 - Offers constant time efficiency (all differentials for an expense similar to one evaluation) for scalar calculations
(A scalar calculation is one where the final result –or output layer, is scalar)

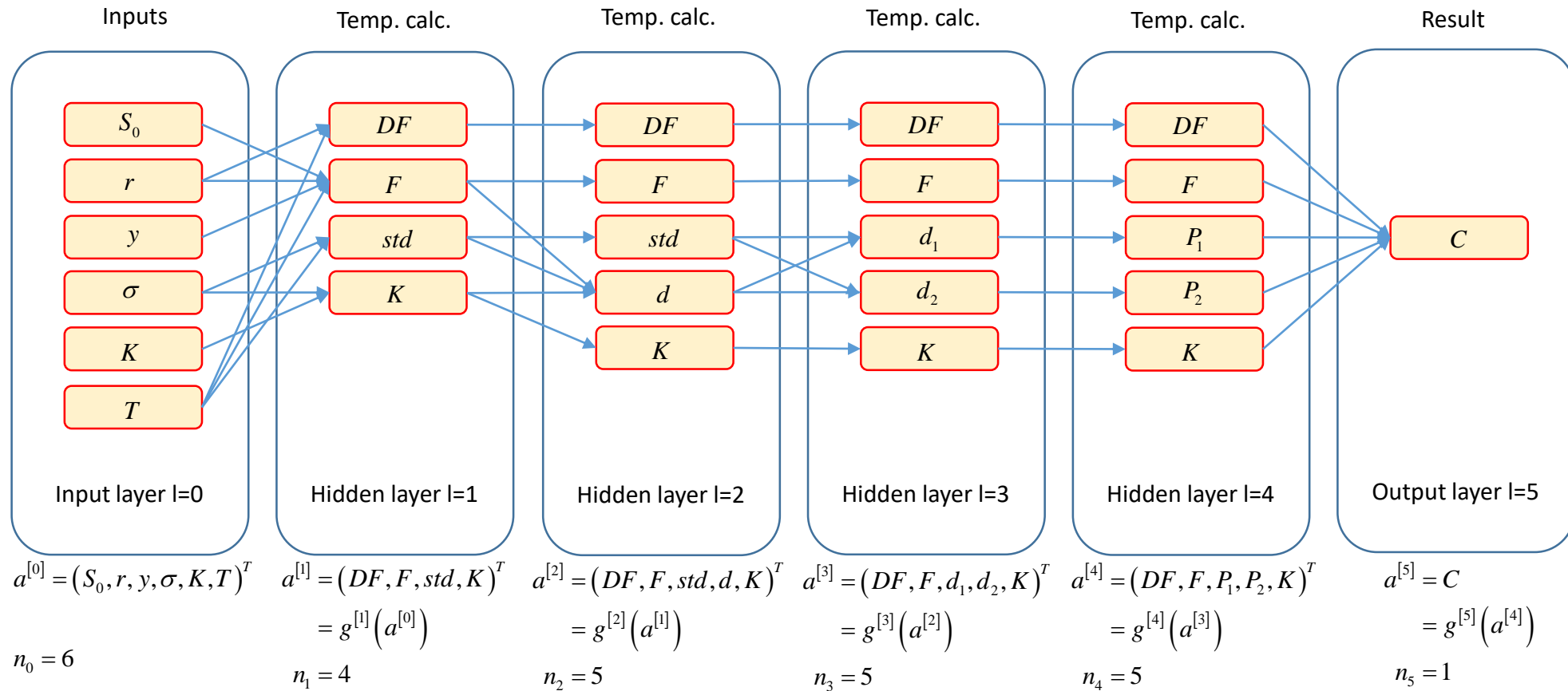


the step functions g are “building blocks” of calculation, which Jacobians are known --- not activations (we are no longer in DL)

Example: Black & Scholes

- Black & Scholes' formula: $C(S_0, r, y, \sigma, K, T) = DF [FN(d_1) - KN(d_2)]$ with
 - Discount factor to maturity: $DF = \exp(-rT)$
 - Forward: $F = S_0 \exp[(r - y)T]$
 - Standard deviation: $std = \sigma\sqrt{T}$
 - Log-moneyness: $d = \frac{\log\left(\frac{F}{K}\right)}{std}$
 - D's: $d_1 = d + \frac{std}{2}$, $d_2 = d - \frac{std}{2}$
 - Probabilities to end in the money, resp. under spot and risk-neutral measures: $P_1 = N(d_1)$, $P_2 = N(d_2)$
 - Call price: $C = DF [FP_1 - KP_2]$

Black & Scholes: evaluation graph



Feed-forward equations

- As in deep nets, for **any** calculation, we have the feed-forward equation defining the evaluation graph:

inputs: $a^{[0]} = x$

feed-forward: $a^{[l]} = g^{[l]}(a^{[l-1]})$

result: $y = a^{[L]}$

- In Black & Scholes, we have the steps:

$$g^{[1]} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \end{pmatrix} = \begin{pmatrix} \exp(-x_2 x_6) \\ x_1 \exp[(x_2 - x_3) x_6] \\ x_4 \sqrt{x_6} \\ x_5 \end{pmatrix} \quad g^{[2]} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \frac{\log(x_2/x_4)}{x_3} \\ x_4 \end{pmatrix} \quad g^{[3]} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \\ x_4 + \frac{x_3}{2} \\ x_4 - \frac{x_3}{2} \\ x_5 \end{pmatrix} \quad g^{[4]} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \\ N(x_3) \\ N(x_4) \\ x_5 \end{pmatrix} \quad g^{[5]} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{pmatrix} = x_1 (x_2 x_3 - x_5 x_4)$$

Differential equations

- For **any** calculation, where we compute the differentials of the result to the inputs: $\frac{\partial y}{\partial x} = \frac{\partial a^{[L]}}{\partial a^{[0]}}$
- If we split the calculation into elementary pieces g , we know the Jacobians $g^{[l]'} = \frac{\partial a^{[l]}}{\partial a^{[l-1]}}$
- In Black & Scholes

$$g^{[1]'} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \end{pmatrix} = \begin{pmatrix} 0 & -x_6 \exp(-x_2 x_6) & 0 & 0 & 0 & -x_2 \exp(-x_2 x_6) \\ \exp[(x_2 - x_3)x_6] & x_1 x_6 \exp[(x_2 - x_3)x_6] & -x_1 x_6 \exp[(x_2 - x_3)x_6] & 0 & 0 & x_1(x_2 - x_3) \exp[(x_2 - x_3)x_6] \\ 0 & 0 & 0 & \sqrt{x_6} & 0 & \frac{x_4}{2\sqrt{x_6}} \\ 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix}^T g^{[2]'} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & \frac{1}{x_2 x_3} & -\frac{\log(x_2/x_4)}{x_3^2} & -\frac{1}{x_3 x_4} \\ 0 & 0 & 0 & 1 \end{pmatrix}^T$$

$$g^{[3]'} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & \frac{1}{2} & 1 & 0 \\ 0 & 0 & -\frac{1}{2} & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix}^T g^{[4]'} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & n(x_3) & 0 & 0 \\ 0 & 0 & 0 & n(x_4) & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix}^T g^{[5]'} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{pmatrix} = ((x_2 x_3 - x_5 x_4), x_1 x_3, x_1 x_2, -x_1 x_5, -x_1 x_4)^T$$

- Note that the Jacobian of the output layer is a vector, as with any scalar calculation

Back-propagation

- We have identified the calculation graph with all its steps g and Jacobians g'
- We have run it once and stored all the intermediate results a
- We can compute all the derivatives of in constant time by back-prop:

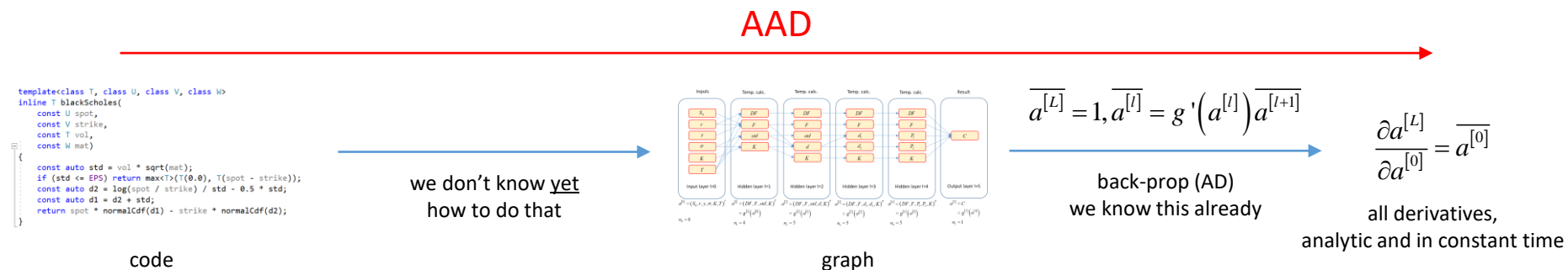
$$\begin{array}{ccc}
 \overline{a^{[L]}} = 1 & \overline{a^{[l]}} = g'(a^{[l]}) \overline{a^{[l+1]}} & \frac{\partial a^{[L]}}{\partial a^{[0]}} = \overline{a^{[0]}} \\
 \swarrow & \uparrow & \searrow \\
 a^{[L]} = \frac{\partial a^{[L]}}{\partial a^{[L]}} = 1 \rightarrow \text{seeding} & & \text{picking results} \\
 & \overline{a^{[l]}} = \underbrace{\frac{\partial a^{[L]}}{\partial a^{[l]}} = \frac{\partial a^{[l+1]}}{\partial a^{[l]}} \frac{\partial a^{[L]}}{\partial a^{[l+1]}}}_{\text{chain rule}} = g'(a^{[l]}) \overline{a^{[l+1]}} \rightarrow \text{backprop} &
 \end{array}$$

- Not that interesting with Black & Scholes: only 6 inputs/derivatives: delta, vega, rho, ...
- The exact same equations apply with large evaluations taking 6,000 inputs --- very interesting then
- Left as an exercise:
Show that Black & Scholes' Jacobians of the previous slide lead to the well known Black-Scholes "Greeks"

Conclusion

- All calculations define a graph
- Like deep nets, all evaluation graphs may be arranged in successive layers, each a an elementary function g of the previous one (which Jacobians are know):

inputs: $a^{[0]} = x$
 feed-forward: $a^{[l]} = g^{[l]}(a^{[l-1]})$
 result: $y = a^{[L]}$
- Once the evaluation graph and its steps g are known, along with their Jacobians g' (and the graph was run once so intermediate results a were stored)
- All derivatives $\frac{\partial a^{[L]}}{\partial a^{[0]}} = \overline{a^{[0]}}$ can be computed in constant time by back-prop $\overline{a^{[L]}} = 1, \overline{a^{[l]}} = g'(a^{[l]}) \overline{a^{[l+1]}}$
- If we can automatically produce a graph from a code, then we can automate the whole process



Recording calculations on tape

Code

- From there on, we work with C++ code
- We only work with simplistic, beginner level C++
- We show that we can still implement AAD with rather spectacular results
- We refer to the Modern Computational Finance book for a professional, generic, parallel, efficient implementation
- **All the code is freely available on GitHub, in the file toyCode.h**

<http://github.com/asavine>

Automatic differentiation

- Differentiation by finite differences
 - Computes derivatives by running evaluation repeatedly
 - Only requires executable evaluation code
 - Linear complexity in the number of differentials
- AAD
 - Computes derivatives by running AD in reverse order over an evaluation graph
 - Requires a graph to traverse – executable code is not enough
 - Constant complexity in the number of differentials : fast differentiation of scalar code
 - To implement AAD we must extract an evaluation graph: sequence of calculations (nodes) and their dependencies (edges)

Building evaluation graphs

- Explicit evaluation graphs
 - Solution implemented in TensorFlow
 - Lets users explicitly build graphs by calling graph creation functions exported to many languages
 - Then ask a TensorFlow session to evaluate or differentiate the graphs efficiently: evaluation and AD, on multicore CPU or GPU
 - Smart and efficient
 - But forces developers to explicitly build graphs in place of calculation code
 - (TensorFlow has a nice API that makes building graphs somewhat similar to coding calculations)
 - But what we want here is take some calculation code and automatically extract its graph
- Source transformation
 - Code that reads and understands evaluation code, then builds graphs and writes backward differentiation code automatically
 - Complex, specialized work similar to writing compilers
 - Variant: template meta-programming and introspection (introduced in chapter 15)

Operator overloading

- Alternative: operator overloading (in languages that support it like C++)
- When code applies operators (like + or *) or math functions (like log or sqrt) to real numbers (*double* type) the corresponding operations are evaluated immediately (or *eagerly*):

```
double x = 1, y = 2;    // x and y are doubles
double z = x + y;       // evaluates x + y and stores result in z
double t = log(x);      // evaluates log(x) and stores result in t
```

- When the same operators are applied to **custom types** (our own type to store real numbers) we decide exactly what happens:

```
class myNumberType
{
    // class definition here
};

myNumberType operator+(const myNumberType& lhs, const myNumberType& rhs)
{
    // this code is executed anytime two numbers of type myNumberType are added
}

myNumberType log(const myNumberType& arg)
{
    // this code is executed anytime log is called on a number of type myNumberType
}

myNumberType x, y;      // x and y are of type myNumberType
myNumberType z = x + y; // executes code in the overloaded operator+()
myNumberType t = log(x); // executes code in the overloaded log()
```

Recording operations

- We apply operator overloading to **record** all operations along with their dependencies:

```
class myNumberType
{
    myNumberType(const double x)
    {
        // constructor initializes value to x and records node
    }
};

myNumberType operator+(const myNumberType& lhs, const myNumberType& rhs)
{
    recordAddition(lhs, rhs); // records addition with dependency on lhs and rhs
}

myNumberType log(const myNumberType& arg)
{
    recordLog(arg); // records log with dependency on arg
}

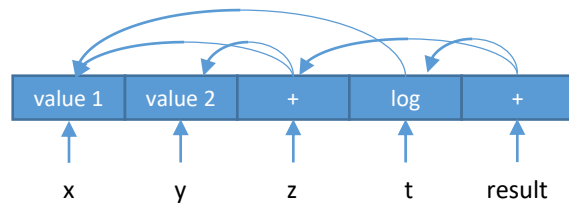
myNumberType x = 1, y = 2; // initializes x to 1 and y to 2 and records them
myNumberType z = x + y; // records addition with dependencies on x and y
myNumberType t = log(x); // records log with dependency on x
myNumberType result = t + z; // records addition with dependency on t and z
```

Lazy evaluation

- When this code is executed:

```
myNumberType x = 1, y = 2;    // initializes x to 1 and y to 2 and records them
myNumberType z = x + y;      // records addition
myNumberType t = log(x);     // records log
myNumberType result = t + z; // records addition
```

- Nothing is calculated, instead the following graph is recorded in memory:



- This sequence can be evaluated later (*lazy evaluation*) or differentiated (applying AD from back to front)
- This is how we build the evaluation graph at run time, by executing calculation code with a custom type for which operators are overloaded to perform recording

Conventional implementation

- We introduce the conventional implementation of AAD
 - Where every math operation: +, -, *, /, pow, log, exp, sqrt, ... is recorded on a data structure called **tape**
 - (It is really an evaluation graph but we call it “tape” in AAD lingo)
- Cutting-edge implementation records whole expressions
 - With template meta-programming and expression templates
 - Resulting in 2x to 5x faster code
 - Sometimes called “tape compression”
 - Extreme instance: “tapeless AD” where we don’t record anything!
 - All implemented in professional code on GitHub and explained in detail in chapter 15
 - In this presentation, we stick with simpler, conventional implementation
- We record every mathematical operation involved in a calculation
 - Every calculation, up to complex Monte-Carlo simulations, is a sequence of +, -, *, /, pow, log, exp, sqrt, ... !
 - We record every single one
 - **Note that all operations have either 0, 1 or 2 arguments**

Simplistic implementation

- In this presentation, we focus on the simplicity of the code
 - We use basic C++ code and disregard efficiency, scalability and best practice
 - Our aim is to explain the key ideas with simplistic code
 - This code works, just not efficiently
- In the book, on the contrary, we build professional, scalable, efficient code in modern C++
- This is particularly important for AAD because of:
 - Recording overhead
 - Every addition, multiplication, etc. produces a record
 - Recording necessarily involves an overhead
 - An efficient implementation must minimize overhead and make recording as efficient as possible
 - Vast memory consumption
 - We store in memory all the operations involved in a large calculation, this is a very large number of records
 - Estimated RAM consumption around 5GB per second
 - Efficient memory and cache management are key to an efficient implementation
- An efficient implementation of conventional AAD is given and explained in chapter 10
 - Effective, custom memory management
 - Recording with minimum overhead, cache efficiency and so on

Record and tape data structures

- A record
 - Stores one operation $y = f(x)$ where $f = +, *, -, /, \log, \sqrt{}, \dots$
 - Knows the number and location of the 0, 1 or 2 arguments x_i
 - Stores the 0, 1 or 2 partial derivatives to its arguments $\partial f / \partial x_i$ so we can apply AD
- The tape stores the sequence of records

```
struct Record
{
    int    numArg;      // number of arguments: 0, 1 or 2
    int    idx1;        // index of first argument on tape
    int    idx2;        // index of second argument on tape
    double der1;        // partial derivative to first argument
    double der2;        // partial derivative to second argument
};

// The tape, declared as a global variable
vector<Record> tape;
```

Custom real number

- Our custom number
 - Holds its value and
 - Knows the index of the corresponding operation on tape
 - May be initialized with a value to create a record (without arguments) on tape

```
struct Number
{
    double value;
    int idx;

    // default constructor does nothing
    Number() {}

    // constructs with a value and record
    Number(const double& x) : value(x)
    {
        // create a new record on tape
        tape.push_back(Record());
        Record& rec = tape.back();

        // reference record on tape
        idx = tape.size() - 1;

        // populate record on tape
        rec.numArg = 0;
    }
};
```

- We overload all mathematical operators and functions to:
 - Evaluate and store result as usual
 - Additionally, record the operation and its derivatives on tape
 - So code is evaluated and recorded at the same time

```
Number operator+(const Number& lhs, const Number& rhs)
{
    // create a new record on tape
    tape.push_back(Record());
    Record& rec = tape.back();

    // compute result
    Number result;
    result.value = lhs.value + rhs.value; // calling double overload

    // reference record on tape
    result.idx = tape.size() - 1;

    // populate record on tape
    rec.numArg = 2;
    rec.idx1 = lhs.idx;
    rec.idx2 = rhs.idx;

    // compute derivatives, both derivatives of addition are 1
    rec.der1 = 1;
    rec.der2 = 1;

    return result;
}
```


Operator overloading

- Similarly, we overload -, *, and /
- Same code exactly, only values and derivatives change

```
Number operator-(const Number& lhs, const Number&rhs)
{
    // ...

    // compute value
    result.value = lhs.value - rhs.value;

    // ...

    // compute derivatives
    rec.der1 = 1;
    rec.der2 = -1;

    // ...
}
```

```
Number operator*(const Number& lhs, const Number&rhs)
{
    // ...

    // compute value
    result.value = lhs.value * rhs.value;

    // ...

    // compute derivatives
    rec.der1 = rhs.value;
    rec.der2 = lhs.value;

    // ...
}
```

```
Number operator/(const Number& lhs, const Number&rhs)
{
    // ...

    // compute value
    result.value = lhs.value / rhs.value;

    // ...

    // compute derivatives
    rec.der1 = 1.0 / rhs.value;
    rec.der2 = - lhs.value / (rhs.value * rhs.value);

    // ...
}
```

On-class operator overloading

- We must also overload +=, -=, *- and /=, as well as unary + and -, on class

The final Number class is therefore:

```
struct Number
{
    double value;
    int idx;

    // default constructor does nothing
    Number() {}

    // constructs with a value and record
    Number(const double& x) : value(x)
    {
        // create a new record on tape
        tape.push_back(Record());
        Record& rec = tape.back();

        // reference record on tape
        idx = tape.size() - 1;

        // populate record on tape
        rec.numArg = 0;
    }

    Number operator +() const { return *this; }
    Number operator -() const { return Number(0.0) - *this; }

    Number& operator +=(const Number& rhs) { *this = *this + rhs; return *this; }
    Number& operator -=(const Number& rhs) { *this = *this - rhs; return *this; }
    Number& operator *=(const Number& rhs) { *this = *this * rhs; return *this; }
    Number& operator /=(const Number& rhs) { *this = *this / rhs; return *this; }
};
```

Function overloading

- Similarly, we overload log, exp, sqrt
- We should really overload all standard math functions
- Code is identical for all functions, only value and derivatives change

```
Number log(const Number& arg)
{
    // create a new record on tape
    tape.push_back(Record());
    Record& rec = tape.back();

    // compute result
    Number result;
    result.value = log(arg.value);

    // reference record on tape
    result.idx = tape.size() - 1;

    // populate record on tape
    rec.numArg = 1;
    rec.idx1 = arg.idx;

    // compute derivative
    rec.der1 = 1.0 / arg.value;

    return result;
}
```

```
Number exp(const Number& arg)
{
    // create a new record on tape
    tape.push_back(Record());
    Record& rec = tape.back();

    // compute result
    Number result;
    result.value = exp(arg.value);

    // reference record on tape
    result.idx = tape.size() - 1;

    // populate record on tape
    rec.numArg = 1;
    rec.idx1 = arg.idx;

    // compute derivative
    rec.der1 = result.value;

    return result;
}
```

```
Number sqrt(const Number& arg)
{
    // create a new record on tape
    tape.push_back(Record());
    Record& rec = tape.back();

    // compute result
    Number result;
    result.value = sqrt(arg.value); // calling double overload

    // reference record on tape
    result.idx = tape.size() - 1;

    // populate record on tape
    rec.numArg = 1;
    rec.idx1 = arg.idx;

    // compute derivative
    rec.der1 = 0.5 / result.value;

    return result;
}
```

Custom function overloading

- Also overload building blocks that are not standard to C++ but frequently applied in applications
- In financial application, we use cumulative normal distributions and normal densities all the time
- The functions are defined in the file gaussians.h in the repo
- We overload here (once again, only change is in value and derivatives):

```
Number normalDens(const Number& arg)
{
    // create a new record on tape
    tape.push_back(Record());
    Record& rec = tape.back();

    // compute result
    Number result;
    result.value = normalDens(arg.value);

    // reference record on tape
    result.idx = tape.size() - 1;

    // populate record on tape
    rec.numArg = 1;
    rec.idx1 = arg.idx;

    // compute derivative
    rec.der1 = - result.value * arg.value;

    return result;
}
```

```
Number normalCdf(const Number& arg)
{
    // create a new record on tape
    tape.push_back(Record());
    Record& rec = tape.back();

    // compute result
    Number result;
    result.value = normalCdf(arg.value); // calling double overload in gaussians.h

    // reference record on tape
    result.idx = tape.size() - 1;

    // populate record on tape
    rec.numArg = 1;
    rec.idx1 = arg.idx;

    // compute derivative
    rec.der1 = normalDens(arg.value);

    return result;
}
```

Avoiding code duplication

- The code for all binary operators and for all unary functions is identical
- Only values and derivatives are different
- This is poor design:
If we change something in the logic, we must consistently modify many different functions
- In the professional code of chapters 10 and 15
we structure code and apply “policy design” (Alexandresku, 2001) to avoid duplication without overhead
- Here, we stick with duplicated code

Comparison operator overloading

- Our custom number must do everything a double does
- In particular, we must be able to compare two Numbers
- Hence, to complete our simple framework, we must also overload comparison operators:

```
bool operator==(const Number& lhs, const Number& rhs) { return lhs.value == rhs.value; }  
bool operator!=(const Number& lhs, const Number& rhs) { return lhs.value != rhs.value; }  
bool operator>(const Number& lhs, const Number& rhs) { return lhs.value > rhs.value; }  
bool operator>=(const Number& lhs, const Number& rhs) { return lhs.value >= rhs.value; }  
bool operator<(const Number& lhs, const Number& rhs) { return lhs.value < rhs.value; }  
bool operator<=(const Number& lhs, const Number& rhs) { return lhs.value <= rhs.value; }
```

Applying the recording framework

- Our simple recording framework is complete, see complete code in the GitHub repo, file toyCode.h
- We may use it to record calculations
- Example: Black and Scholes

```
inline double blackScholes(  
    // input layer 0  
    const double spot, const double rate, const double yield, const double vol, const double strike, const double mat)  
{  
    /* layer 1 */      double df = exp(-rate * mat), fwd = spot * exp((rate - yield) * mat), std = vol * sqrt(mat);  
    /* layer 2 */      double d = log(fwd / strike) / std;  
    /* layer 3 */      double d1 = d + 0.5 * std, d2 = d - 0.5 * std;  
    /* layer 4 */      double p1 = normalCdf(d1), p2 = normalCdf(d2);  
    /* output layer 5 */ return df * (fwd * p1 - strike * p2);  
}
```

Instrumenting computation code

- To record a Black & Scholes calculation, we must call it with our number type
- This is called *instrumentation*
- We can replace all doubles by Numbers:

```
inline Number blackScholes(  
    // input layer 0  
    const Number spot, const Number rate, const Number yield, const Number vol, const Number strike, const Number mat)  
{  
    /* layer 1 */      Number df = exp(-rate * mat), fwd = spot * exp((rate - yield) * mat), std = vol * sqrt(mat);  
    /* layer 2 */      Number d = log(fwd / strike) / std;  
    /* layer 3 */      Number d1 = d + 0.5 * std, d2 = d - 0.5 * std;  
    /* layer 4 */      Number p1 = normalCdf(d1), p2 = normalCdf(d2);  
    /* output layer 5 */ return df * (fwd * p1 - strike * p2);  
}
```


Instrumenting computation code

- Better solution: template code on number representation type

```
template <class T> inline T blackScholes(  
    // input layer 0  
    const T spot, const T rate, const T yield, const T vol, const T strike, const T mat)  
{  
    /* layer 1 */      T df = exp(-rate * mat), fwd = spot * exp((rate - yield) * mat), std = vol * sqrt(mat);  
    /* layer 2 */      T d = log(fwd / strike) / std;  
    /* layer 3 */      T d1 = d + 0.5 * std, d2 = d - 0.5 * std;  
    /* layer 4 */      T p1 = normalCdf(d1), p2 = normalCdf(d2);  
    /* output layer 5 */ return df * (fwd * p1 - strike * p2);  
}
```

- Best practice: produce templated code in the first place

- To evaluate only, call with doubles as arguments

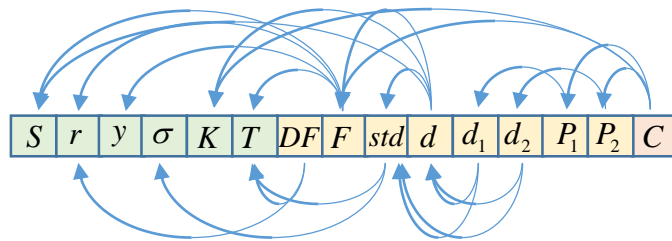
```
double spot = 100, rate = 0.02, yield = 0.05, vol = 0.2, strike = 110, mat = 2; // initializes inputs  
auto result = blackScholes(spot, rate, yield, vol, strike, mat);                // evaluates operations
```

- To evaluate and record code, call with Numbers as arguments

```
Number spot = 100, rate = 0.02, yield = 0.05, vol = 0.2, strike = 110, mat = 2; // initializes and records inputs  
auto result = blackScholes(spot, rate, yield, vol, strike, mat);                // evaluates and records operations
```

Tape vs graph

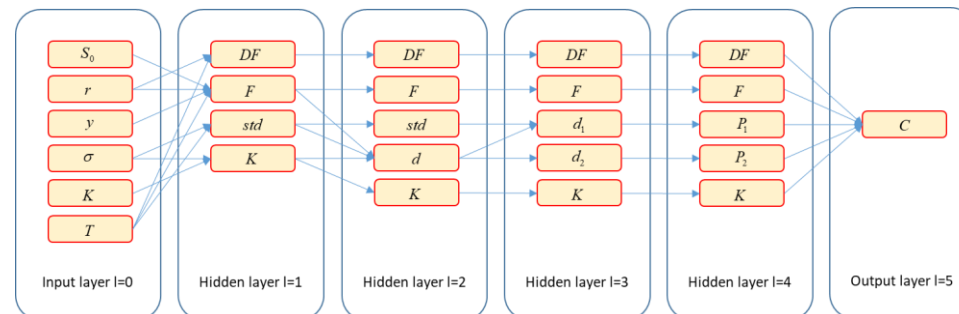
- Operator overloading gives us a tape of the mathematical operations involved in a calculation
- After execution of the templated code for Black & Scholes, we get the following tape:



Where we recall that operator overloading also gave us:

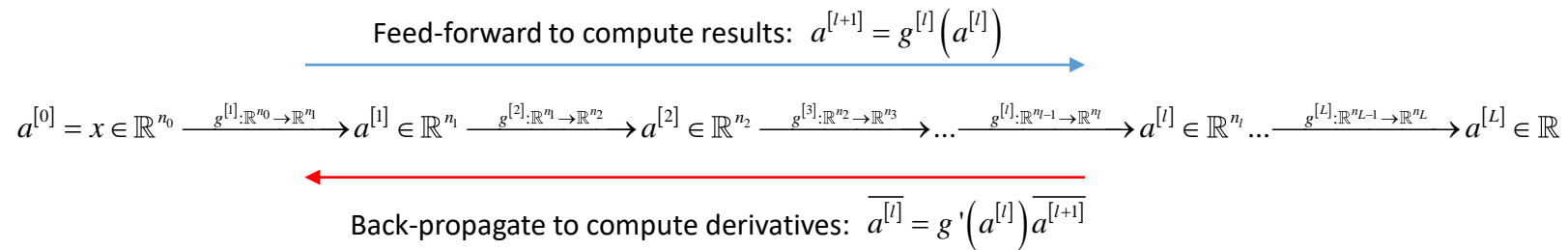
- The operations from ancestors A (arguments) to successors S (results)
- And their derivatives $\frac{\partial S}{\partial A}$

- A tape is not exactly the same as an evaluation graph, where operations are neatly arranged by layer:



Tape vs graph (2)

- The tape and the graph however express the exact same information, differently
- We can always turn a tape into a graph or vice-versa
- If we had a graph, we could run back-prop to compute all derivatives in constant time



- We could turn our tape into a graph, but this would be a waste of precious computational resources
- Instead, we work with a slightly different version of the back-prop/AD algorithm
- So we can back-propagate adjoints directly through the tape (although it is not arranged in layers)
- Using the notions of *ancestor* (argument to an operation) and *successor* (result of an operation)
- AAD is the sum of a *taping framework* (e.g. operator overloading) and back-propagation through the tape

AAD

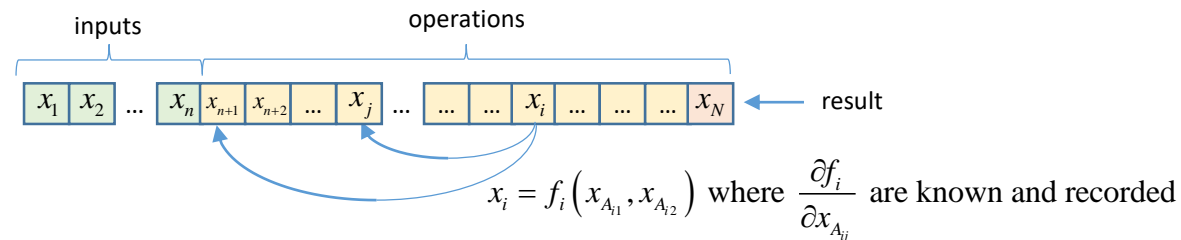
Ancestors

- We call the instrumented instance of our code:

```
Number spot = 100, rate = 0.02, yield = 0.05, vol = 0.2, strike = 110, mat = 2; // initializes and records inputs
auto result = blackScholes(spot, rate, yield, vol, strike, mat); // evaluates and records operations
cout << result.value; // 5.03705
```

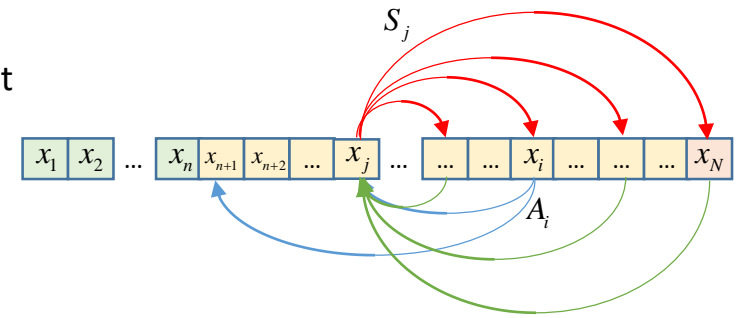
- Which evaluates the calculation and records all operations on tape
 - Denote f_i the operation number i , this is a function of 0, 1 or 2 arguments
 - The arguments must also be on tape, with indices $< i$
 - Denote A_i the set of indices of the ancestors (arguments) to f_i - this is a set of 0, 1 or 2 indices, all $< i$
 - Denote n the number of inputs and N the number of operations on tape, including inputs
 - Denote x_i the result of operation i and $y = x_N$ the final result
- Note that all the local derivatives $\frac{\partial f_i}{\partial x_j}$ for all $j \in A_i$ have been computed and recorded on tape, on evaluation, by our operators

- Our tape therefore looks like:



Successors

- Successors: j is a successor to i if i is an ancestor to j
 - Denote S_j the set of indices of the *successors* of x_j on tape
 - These are the indices i of all functions f_i , calculated **after** x_j that use x_j as an argument
 - Formally: $S_j = \{i > j, j \in A_i\}$
 - Note that S_j may contain many indices, although the size of A_i is 2 or less
 - Whereas an empty S_j reveals an unused input or intermediate result



- Adjoint equation
 - Denote $\bar{x}_i \equiv \frac{\partial y}{\partial x_i} = \frac{\partial x_N}{\partial x_i}$ the adjoint of x_i
 - Then (evidently): $\bar{x}_N = 1$
 - And in a direct application of the chain rule: $\bar{x}_j = \sum_{i \in S_j} \frac{\partial f_i}{\partial x_j} \bar{x}_i$ because $\frac{\partial y}{\partial x_j} = \sum_{i \in S_j} \frac{\partial y}{\partial x_i} \frac{\partial x_i}{\partial x_j}$ and we recall that all the $\frac{\partial f_i}{\partial x_j}$ are all on tape

Adjoint accumulation V1

- Adjoint accumulation therefore satisfy a backward recursion

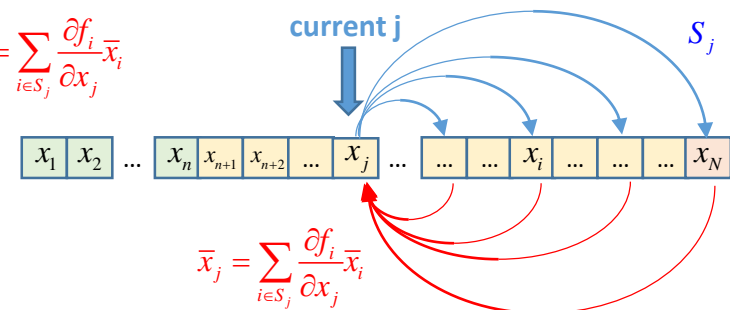
- The last adjoint $\bar{x}_N = 1$ is given

- All other adjoints are a weighted sum of future adjoints: for all j
$$\bar{x}_j = \sum_{i \in S_j} \frac{\partial f_i}{\partial x_j} \bar{x}_i$$

- Hence, the following (flawed) algorithm:

Starting with $\bar{x}_N = 1$, compute for every j , in the reverse order, from $N-1$ to 1 :
$$\bar{x}_j = \sum_{i \in S_j} \frac{\partial f_i}{\partial x_j} \bar{x}_i$$

Summing over the successors of x_j :

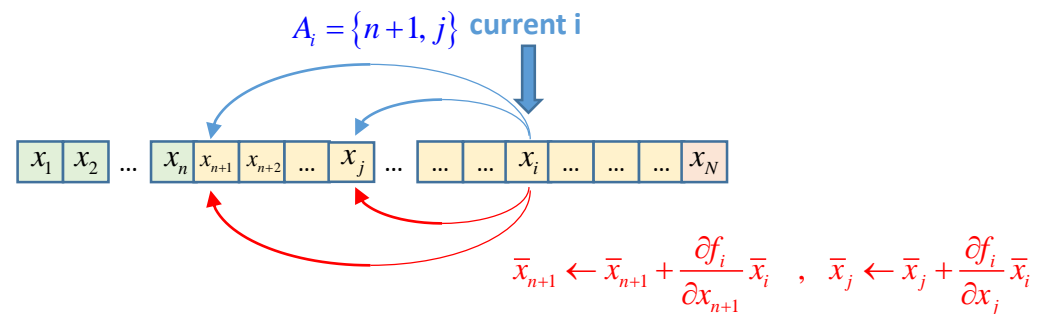


- Problems:

- Complexity: a node may have many successors to sum over
- Impracticality: ancestors don't store their successors, it is the successors who store their (up to 2) ancestors and partial derivatives

Adjoint accumulation V2

- A more efficient/practical means of computing the same numbers:
 - The last adjoint $\bar{x}_N = 1$ is still given
 - Initialize all other adjoints to 0
 - Traverse tape in the reverse order, as previously, accumulating the sum $\bar{x}_j = \sum_{i \in S_j} \frac{\partial f_i}{\partial x_j} \bar{x}_i$ **from successors**



- Problems resolved:
 - Max complexity 2: a node may have up to 2 ancestors
 - Successors who store their ancestors and partial derivatives

Adjoint accumulation algorithm

1. Initialize all adjoints to 0 and the last adjoint to 1 (this is called *seeding* the tape): $\bar{x}_j = \delta_{N-j}$
 2. Repeat for i iterating backwards from N to 0 : for all $j \in A_i$: $\bar{x}_j \leftarrow \bar{x}_j + \frac{\partial f_i}{\partial x_j} \bar{x}_i$
 3. The differentials of the calculation to its inputs x_1, x_2, \dots, x_n are, by definition, $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n$
- This algorithm is called (reverse) adjoint propagation
 - AAD is the sum of a recording framework and an implementation of adjoint propagation

Adjoint accumulation code

```
vector<double> calculateAdjoints(Number& result)
{
    // initialization
    vector<double> adjoints(tape.size(), 0.0); // initialize all to 0
    int N = result.idx;                       // find N
    adjoints[N] = 1.0;                         // seed aN = 1

    // backward propagation
    for(int j=N; j>0; --j) // iterate backwards over tape
    {
        if (tape[j].numArg > 0)
        {
            adjoints[tape[j].idx1] += adjoints[j] * tape[j].der1; // propagate first argument

            if (tape[j].numArg > 1)
            {
                adjoints[tape[j].idx2] += adjoints[j] * tape[j].der2; // propagate second argument
            }
        }
    }

    return adjoints;
}
```

Application to Black & Scholes

- After we record a tape by calling the instrumented instance of our code:

```
Number spot = 100, rate = 0.02, yield = 0.05, vol = 0.2, strike = 110, mat = 2; // initializes and records inputs
auto result = blackScholes(spot, rate, yield, vol, strike, mat); // evaluates and records operations
cout << result.value; // 5.03705
```

- We proceed with adjoint propagation:

```
// propagate adjoints
vector<double> adjoints = calculateAdjoints(result);
```

- The code works nicely
and produces correct values

```
// show derivatives
cout << "Derivative to spot (delta) = " << adjoints[spot.idx] << endl; // 0.309
cout << "Derivative to rate (rho) = " << adjoints[rate.idx] << endl; // 51.772
cout << "Derivative to dividend yield = " << adjoints[yield.idx] << endl; // -61.846
cout << "Derivative to volatility (vega) = " << adjoints[vol.idx] << endl; // 46.980
cout << "Derivative to strike (-digital) = " << adjoints[strike.idx] << endl; // -0.235
cout << "Derivative to maturity (-theta) = " << adjoints[mat.idx] << endl; // 1.321
```

Complexity

- One *evaluation* of the calculation
 - Sweeps forward through the sequence of its operations
 - Executes every operation exactly once
 - Therefore its complexity is N , the number of records that end up on tape
- Adjoint propagation
 - Also sweeps through the sequence of operations, backward
 - Executes 0, 1 or 2 operations on every record, depending on the number of arguments
 - Therefore, as advertised, AAD computes **all n** differentials in **constant time**
 - In theory, all differentials are propagated in less than 2x one evaluation
 - In addition, the calculation must be evaluated (and recorded) first so the theoretical upper bound is 3x one evaluation
 - Due to recording and tape traversal overhead, a good implementation generally produces many differentials in 4x to 10x
- Our professional code from chapters 10, 12 and 15 beats the theoretical bound!
 - Recall from the demonstration, one evaluation = 0.8sec, 1,081 differentials = 1.5sec, less than 2x one evaluation
 - Due to “selective instrumentation”, a strong optimization, explained, along many others, in chapter 12

Conclusion

- We implemented AAD in the simplest possible manner, scratching the surface of possibilities
- Part III (Chapters 8 to 15) gives the details of a complete, professional, efficient implementation
 - How to minimize recording overhead
 - Efficient memory management constructs
 - Apply check-pointed AAD to differentiate a calculation piece by piece to mitigate RAM footprint and cache inefficiency
 - Efficiently differentiate non-scalar calculations that return multiple results
 - Cutting-edge implementation with template meta-programming and expression templates, faster by 2x to 5x
 - Parallel implementation
 - Advise for debugging and optimization
 - And much more
- Still the simplistic code works and produces the correct values
- Not very interesting for Black & Scholes
 - Fast, analytic evaluation
 - Only 6 differentials to compute
- Next, we apply the framework to a barrier option in Dupire Monte-Carlo
 - Long, complex evaluation
 - 1,081 differentials to compute

AAD for financial simulations

Simple simulation code

- We implement a simplistic simulation code for a barrier option in Dupire's model
 - Recall local volatility is given in a matrix and bi-linearly interpolated in spot and time
- We need the following pieces, which we assume are given here
 - A matrix class to hold local volatilities (matrix.h in the repo, chapters 1 and 2 in the book)
 - A bi-linear interpolation function (interp.h in the repo, chapter 6, section 6.4 in the book)
 - Random number generators to produce independent Gaussian increments (chapters 5 and 6)
- The code is templated on the real number representation type

Simulation code, version 1

// Signature

```
template <class T>
inline T toyDupireBarrierMc(
    // Spot
    const T S0,
    // Local volatility
    const vector<T> spots,
    const vector<T> times,
    const matrix<T> vols,
    // Product parameters
    const T maturity,
    const T strike,
    const T barrier,
    // Number of paths and time steps
    const int Np,
    const int Nt,
    // Initialized random number generator
    RNG& random)
```

// Implementation

```
// Initialize
T result = 0;
// double because the RNG is not templated (and doesn't need to be, see chapter 12)
vector<double> gaussianIncrements(Nt);
const T dt = maturity / Nt, sdt = sqrt(dt);

// Loop over paths
for (int i = 0; i < Np; ++i)
{
    // Generate Nt Gaussian Numbers
    random.nextG(gaussianIncrements);
    // Euler's scheme, step by step
    T spot = S0, time = 0;
    bool alive = true;
    for (size_t j = 0; j < Nt; ++j)
    {
        // Interpolate volatility
        const T vol = interp2D(spots, times, vols, spot, time);
        time += dt;
        // Simulate return
        spot *= exp(-0.5 * vol * vol * dt + vol * sdt * gaussianIncrements[j]);
        // Monitor barrier
        if (spot > barrier)
        {
            alive = false;
            break;
        }
    }
    // Payoff
    if (alive && spot > strike) result += spot - strike;
} // paths

return result / Np;
```


A simplistic code

- This code “does the job” but is not acceptable by professional standards
- The code is specific to Dupire’s model and an up & out call, therefore not scalable
 - To price another product (Asian option, Ratchet option, ...) copy the code and change the lines that evaluate payoffs
 - To price in another model (Heston, SLV, ...) copy the code and change the lines that generate the path
 - End up with many different functions implementing the same simulation logic for different couples of models and products
 - To modify the simulation logic, consistently change all the functions! This is obviously not viable
- Chapter 6 teaches a professional architecture for generic simulation libraries
 - Encapsulate scenario generation in Model objects
 - Encapsulate payoff evaluation in Product objects
 - Encapsulate simulation logic in a generic Monte-Carlo engine
 - Code every model and every product exactly once, mix and match at run time
- We stick with the simplistic code for demonstration purposes

An inefficient code

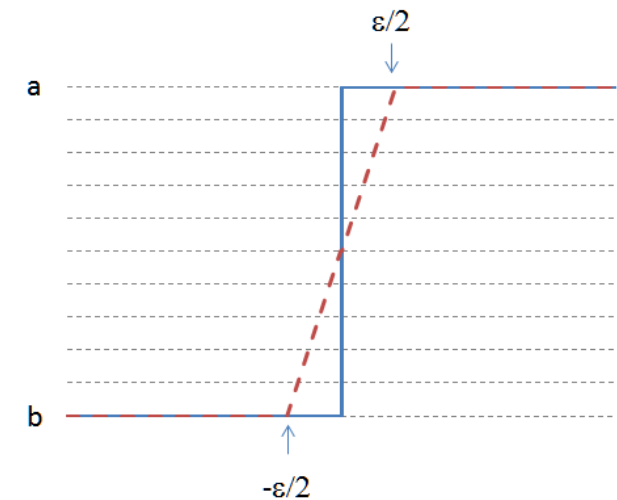
- The code does too much work repeatedly during simulations
 - Key to efficient Monte-Carlo code:
Do as much work as possible once, on initialization, and as little as possible repeatedly, during simulations
 - Example: we perform an expensive bi-linear interpolation in the innermost loop, for every path, on every time step
 - We interpolate in spot and time, spot is stochastic (scenario dependent), time is not
 - Therefore we can (and should) pre-interpolate in time on initialization
And perform only 1D interpolations in the innermost loop
- Chapter 6 teaches and builds fully optimized code
- The code is serial, it executes sequentially on one core
 - Since even our phones are multi-core today, professional code must be parallel
 - With Monte-Carlo simulations, it is relatively easy to obtain a speed-up by the number of physical cores
- Chapter 3 teaches modern parallel C++, chapter 7 builds a professional *parallel* simulation library
And section 12.5 instruments it with AAD in parallel
- For the purpose of demonstration, we stick with the not so efficient, serial version

Smoothing barrier options

- Discretely monitored barrier options are discontinuous, their value jumps to 0 at the barrier
 - Therefore, our code is not differentiable
 - AAD does not help: it cannot perform the impossible task of differentiating a discontinuous function
 - With finite differences, barrier risks are unstable, with AAD they are all zero
 - Find the reason why as an exercise!
- Therefore traders always smooth discontinuous transactions (barriers, digitals etc.)
- Smoothing = applying a close, continuous approximation in place of the discontinuous function
- Smoothing in finance, and its connection to fuzzy logic, are explained in the presentation:
<http://www.slideshare.net/AntoineSavine/stabilise-risks-of-discontinuous-payoffs-with-fuzzy-logic>
freely available on slideShare
- Here, we briefly explain the “smooth barrier” algorithm, usually applied on derivatives desks

Smooth barrier

- Hard barrier
 - 100% dead above the barrier, 100% alive below the barrier
 - Hence, discontinuous
- Soft barrier
 - 100% dead above the barrier **plus epsilon**, 100% alive below the barrier **minus epsilon**
 - In between, lose a fraction of notional interpolated between (barrier-epsilon,0) and (barrier+epsilon,1)
 - And continue with the remaining notional
 - Hence, continuous
- Smoothing and fuzzy logic
 - Like Schrodinger's cat, the transaction is in a superposition of dead and alive states
 - Smoothing is achieved by replacing sharp logic (dead or alive?) by fuzzy logic (how much alive?)
 - More in the presentation



Simulation code with smooth barrier

```
template <class T>
inline T toyDupireBarrierMc(
    // Spot
    const T S0,
    // Local volatility
    const vector<T> spots,
    const vector<T> times,
    const matrix<T> vols,
    // Product parameters
    const T maturity,
    const T strike,
    const T barrier,
    // Number of paths and time steps
    const int Np,
    const int Nt,
    // Smoothing
    const T epsilon,
    // initialized random number generator
    RNG& random)

// Initialize
T result = 0;
// double because the RNG is not templated (and doesn't need to be, see chapter 12)
vector<double> gaussianIncrements(Nt);
const T dt = maturity / Nt, sdt = sqrt(dt);

// Loop over paths
for (int i = 0; i < Np; ++i)
{
    // Generate Nt Gaussian Numbers
    random.nextG(gaussianIncrements);
    // Step by step
    T spot = S0, time = 0;
    /* bool alive = true; */ T alive = 1.0; // alive is a real number in (0,1)
    for (size_t j = 0; j < Nt; ++j)
    {
        // Interpolate volatility
        const T vol = interp2D(spots, times, vols, spot, time);
        time += dt;
        // Simulate return
        spot *= exp(-0.5 * vol * vol * dt + vol * sdt * gaussianIncrements[j]);
        // Monitor barrier
        /* if (spot > barrier) { alive = false; break; } */
        if (spot > barrier + epsilon) { alive = 0.0; break; } // definitely dead
        else if (spot < barrier - epsilon) { /* do nothing */ }; // definitely alive
        else /* in between, interpolate */ alive *= 1.0 - (spot - barrier + epsilon) / (2 * epsilon);
    }
    // Payoff paid on surviving notional
    /* if (alive && spot > strike) result += spot - strike; */ if (spot > strike) result += alive * (spot - strike);
} // paths
return result / Np;
```

Simulation code

- The toy simulation code is found on the file ToyCode.h in the repo
- To be compared with professional code in the files with names prefixed by “mc”
- Our simple code returns the same result as the professional code
- It is twice slower than the serial version of the professional code
- On a quad-core computer, it is 8x slower than the parallel version
- Next, we differentiate it with our simple AAD framework

Differentiation

Just like we did for Black & Scholes, to compute differentials, we:

1. Initialize the inputs as Numbers
Which also records them on tape
2. Call our templated evaluation code, instantiated with the Number type
Which performs the evaluation *and* records operations on tape
3. Propagate adjoints backwards through the tape
4. Pick differentials as the adjoints of the parameters

Differentiation code

```
void toyDupireBarrierMcRisks(
    const double S0, const vector<double> spots, const vector<double> times, const matrix<double> vols,
    const double maturity, const double strike, const double barrier,
    const int Np, const int Nt, const double epsilon, RNG& random,
    /* results: value and dV/dS, dV/d(local vols) */ double& price, double& delta, matrix<double>& vegas)
{
    // 1. Initialize inputs

    Number nS0(S0), nMaturity(maturity), nStrike(strike), nBarrier(barrier), nEpsilon(epsilon);
    vector<Number> nSpots(spots.size()), nTimes(times.size());
    matrix<Number> nVols(vols.rows(), vols.cols());

    for (int i = 0; i < spots.size(); ++i) nSpots[i] = Number(spots[i]);
    for (int i = 0; i < times.size(); ++i) nTimes[i] = Number(times[i]);
    for (int i = 0; i < vols.rows(); ++i) for (int j = 0; j < vols.cols(); ++j) nVols[i][j] = Number(vols[i][j]);

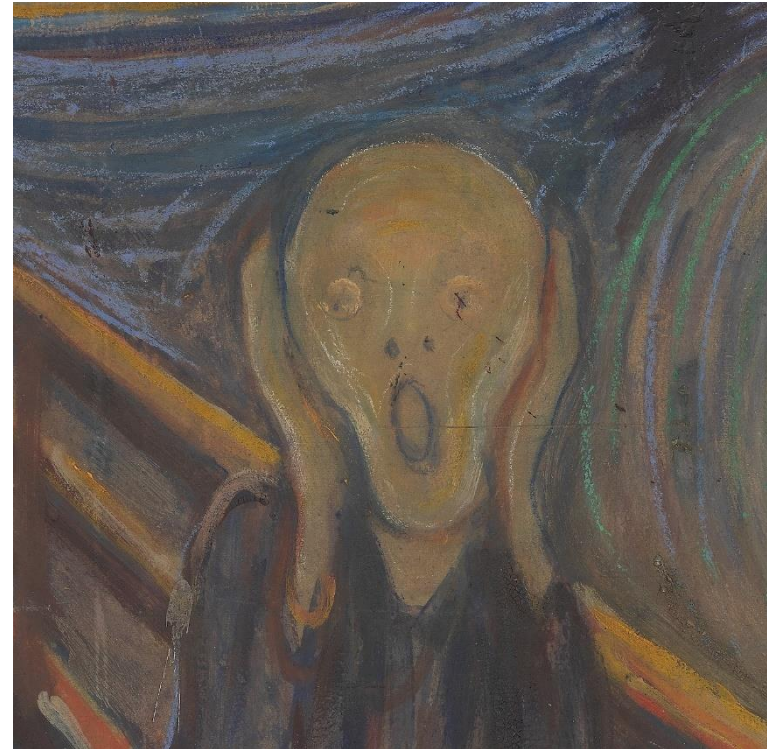
    // 2. Call instrumented evaluation code, which evaluates the barrier option price and records all operations
    Number nPrice = toyDupireBarrierMc(nS0, nSpots, nTimes, nVols, nMaturity, nStrike, nBarrier, Np, Nt, nEpsilon, random);

    // 3. Adjoint propagation, the exact same code as before, should be encapsulated in a dedicated function
    vector<double> adjoints = calculateAdjoints(nPrice);

    // 4. Pick results
    price = nPrice.value;
    delta = adjoints[nS0.idx];
    for (int i = 0; i < vols.rows(); ++i) for (int j = 0; j < vols.cols(); ++j) vegas[i][j] = adjoints[nVols[i][j].idx];
}
```


Testing the code

- We run the code in the same context as the initial demonstration but with 100,000 paths instead of 500,000
- The computer runs out of memory and crashes!
- Running AAD on a simulation with 100,000 paths consumes an insane amount of RAM
- Even on a computer with enough memory, such large tape is cache inefficient



Solution in principle

- Run a series of risks on mini-batches of say, 1024 paths and average in the end
- Wipe the tape in between mini-batches
- The average of differentials is the differential of the average
- So we get the same results while reducing memory footprint to operations recorded over 1,024 paths
- Note with mini-batches of size 1, this is known as “path-wise differentiation”
- This is also a particular, and simple case of the general check-pointing algorithm explained in chapter 13

Solution in code

- Rename our function DupireRisksMiniBatch(), call it sequentially from a wrapper function

```
void toyDupireBarrierMcRisks(
    const double S0, const vector<double> spots, const vector<double> times, const matrix<double> vols,
    const double maturity, const double strike, const double barrier,
    const int Np, const int Nt, const double epsilon, RNG& random,
    /* results: value and dV/dS, dV/d(local vols) */ double& price, double& delta, matrix<double>& vegas)
{
    price = delta = 0;
    for (int i = 0; i < vegas.rows(); ++i) for (int j = 0; j < vegas.cols(); ++j) vegas[i][j] = 0;
    double batchPrice, batchDelta; matrix<double> batchVegas(vegas.rows(), vegas.cols());
    int pathsToGo = Np, pathsPerBatch = 1024;
    // calculate batch sensitivities sequentially
    while (pathsToGo > 0)
    {
        // wipe tape
        tape.clear();

        // do mini batch
        int paths = min(pathsToGo, pathsPerBatch);
        dupireRisksMiniBatch(S0, spots, times, vols, maturity, strike, barrier, paths, Nt, epsilon, random, batchPrice, batchDelta, batchVegas);

        // update results
        price += batchPrice * paths / Np;
        delta += batchDelta * paths / Np;
        for (int i = 0; i < vegas.rows(); ++i) for (int j = 0; j < vegas.cols(); ++j) vegas[i][j] += batchVegas[i][j] * paths / Np;

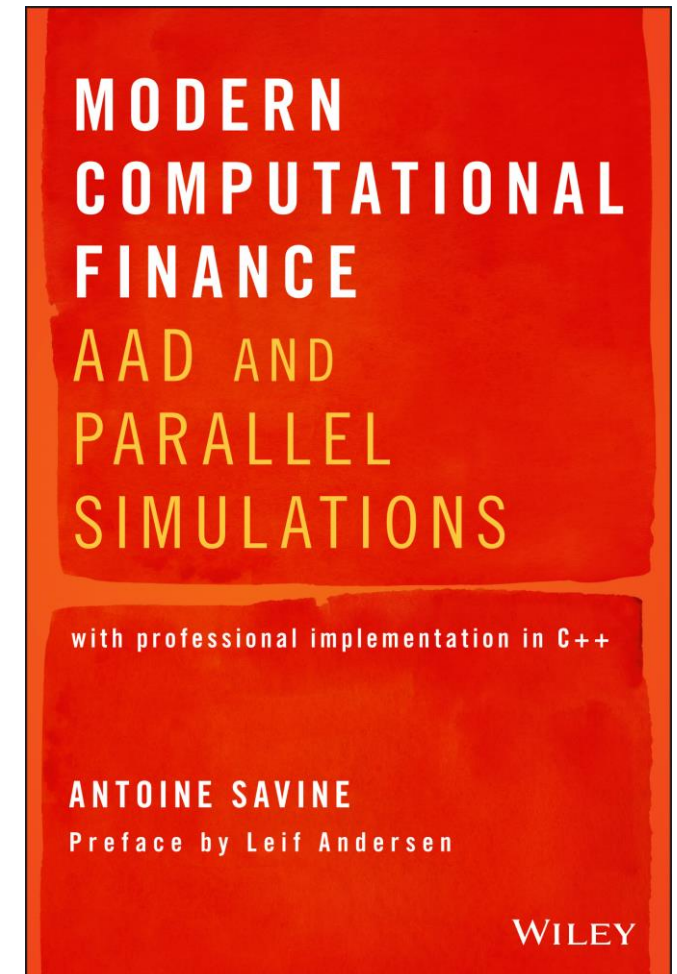
        pathsToGo -= paths;
    }
}
```

Performance

- With 100,000 paths, 156 steps, we compute the 1,081 differentials in around 7 seconds
- This is 1,081 differentials in the time of around 6 evaluations
- This is a very remarkable result, especially with such simplistic code
- Try it yourself with the toy code in the repo!
- This being said, the professional code is around 8 times faster in serial mode, 32 times faster in parallel mode on a quad-core laptop

Conclusion

- We learned AAD in principle and in code
And applied it to machine learning and finance
- But we really just scratched the surface
- For example:
 - How to efficiently differentiate the multiple results of non-scalar functions?
 - How to compute risks not on model parameters like local volatilities
But on tradable market variables like implied volatilities?
 - How to implement AAD in modern C++ and manage memory efficiently?
 - How to implement AAD over parallel simulations and run it at least 32x faster?
- The answers, and much more, are in the book



Thank you for your attention. Find the slides here:

go to <http://github.com/asavine>

click 'CompFinance'

download slides 'Intro2AADinMachineLearningAndFinance.pdf'

download examples in the 'Workshop' folder

see tutorial in the 'xlCpp' folder for exporting C++ code to Excel

see basic AAD code in C++ in the file 'toyCode.h'

the rest of the repo is companion code for the book [Modern Computational Finance](#)

follow asavine on GitHub and watch the CompFinance repo to be notified of updates