

Automatic Differentiation (2)

Slides Prepared By:

Atılım Güneş Baydin
gunes@robots.ox.ac.uk

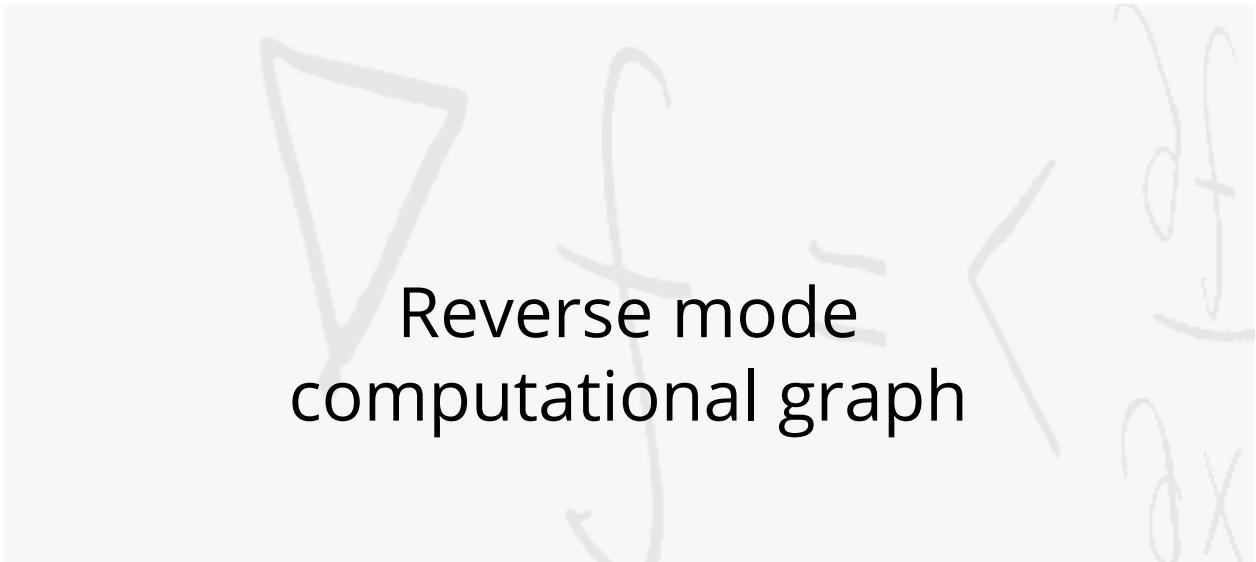
Summary of last lecture

- Derivatives in machine learning
- Review of essential concepts
 - derivative, partial derivative, total derivative, gradient, Jacobian, matrix calculus, etc.
- How do we compute derivatives
 - Manual, symbolic, numerical
- Automatic differentiation
- Computational graphs and propagation

Today

- The reverse mode (backprop) computational graph
 - What gets propagated?
- Implementation
 - Where does the graph come from?
 - Strategies and performance tips
- Advanced concepts
 - Nesting, higher-order derivatives
 - Checkpointing

3



Reverse mode
computational graph

4

Normal PDF

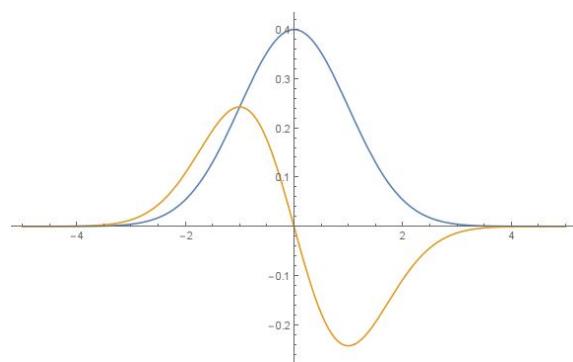
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

5

Normal PDF

$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

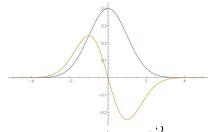
$$\frac{\partial f}{\partial x} = \frac{(\mu - x)e^{-\frac{(\mu-x)^2}{2\sigma^2}}}{\sqrt{2\pi}\sigma^3} \quad \frac{\partial f}{\partial \mu} = \frac{(x - \mu)e^{-\frac{(\mu-x)^2}{2\sigma^2}}}{\sqrt{2\pi}\sigma^3} \quad \frac{\partial f}{\partial \sigma} = -\frac{(\sigma - x + \mu)(\sigma + x - \mu)e^{-\frac{(\mu-x)^2}{2\sigma^2}}}{\sqrt{2\pi}\sigma^4}$$



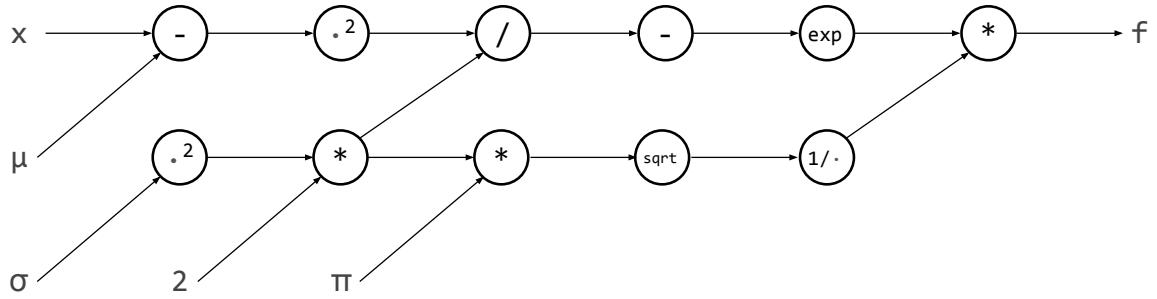
6

Normal PDF

$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



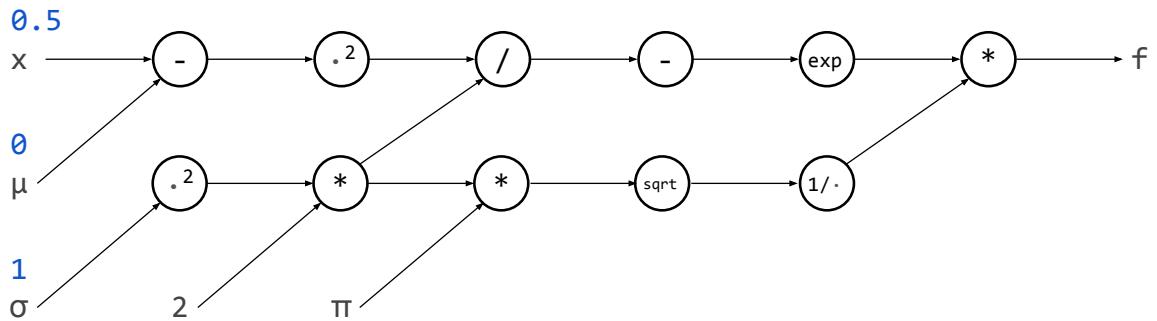
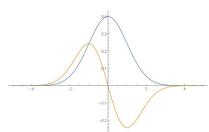
$$\frac{\partial f}{\partial x} = \frac{(\mu - x)e^{-\frac{(\mu-x)^2}{2\sigma^2}}}{\sqrt{2\pi}\sigma^3} \quad \frac{\partial f}{\partial \mu} = \frac{(x - \mu)e^{-\frac{(x-\mu)^2}{2\sigma^2}}}{\sqrt{2\pi}\sigma^3} \quad \frac{\partial f}{\partial \sigma} = -\frac{(\sigma - x + \mu)(\sigma + x - \mu)e^{-\frac{(x-\mu)^2}{2\sigma^2}}}{\sqrt{2\pi}\sigma^4}$$



7

Normal PDF

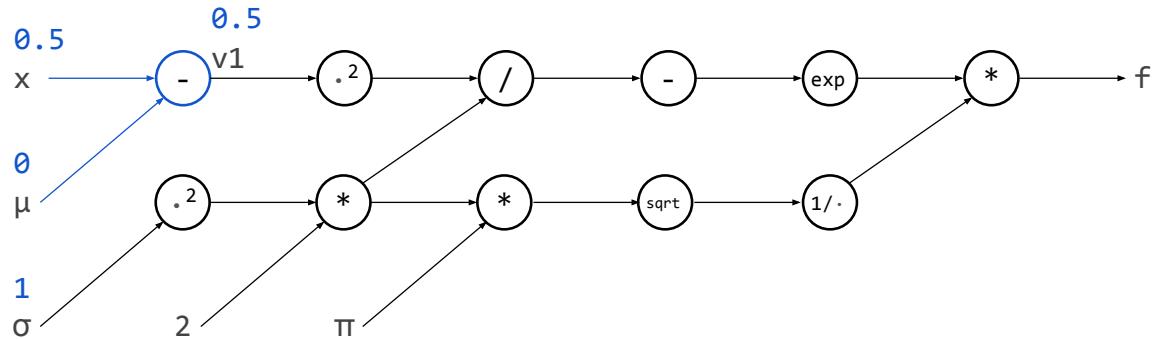
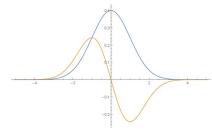
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



8

Normal PDF

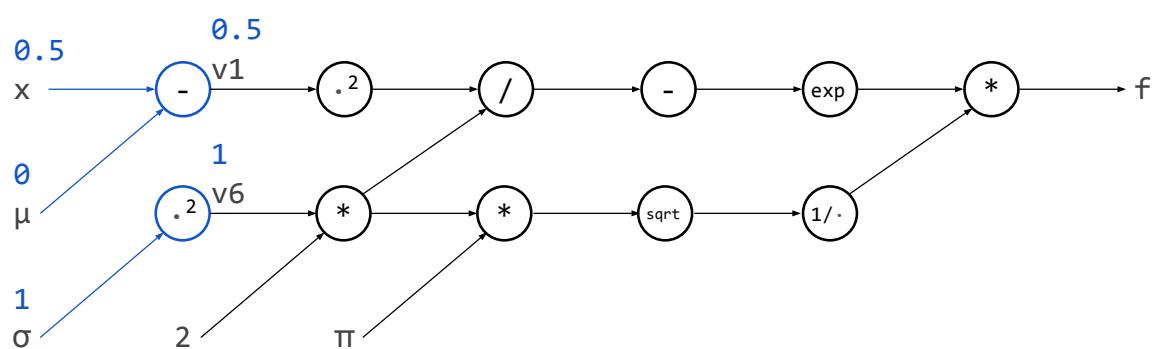
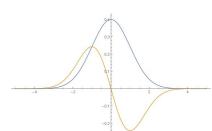
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



9

Normal PDF

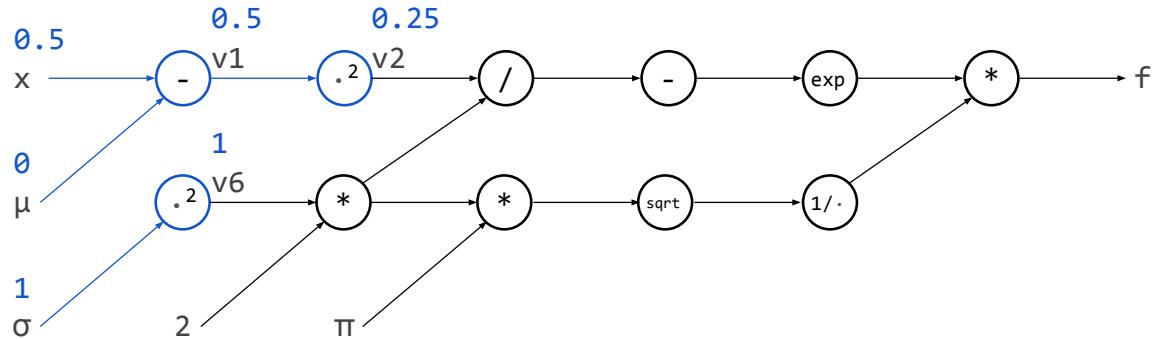
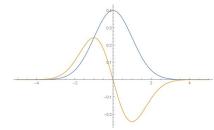
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



10

Normal PDF

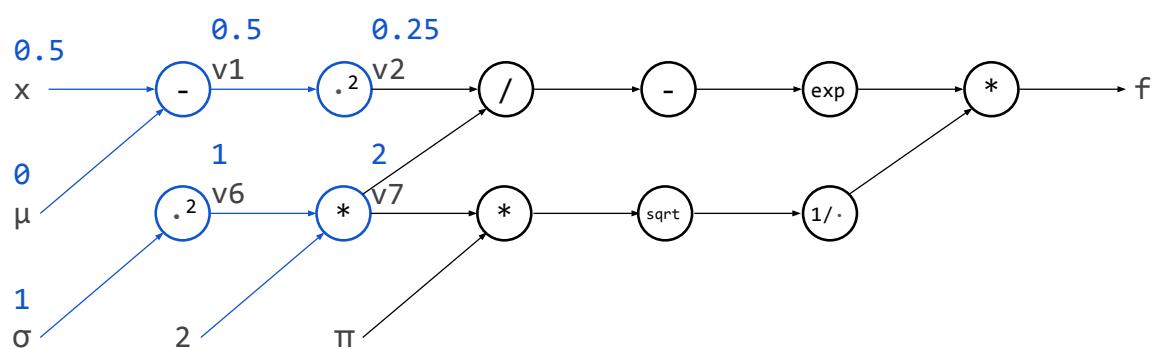
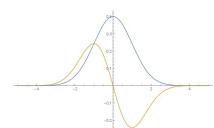
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



11

Normal PDF

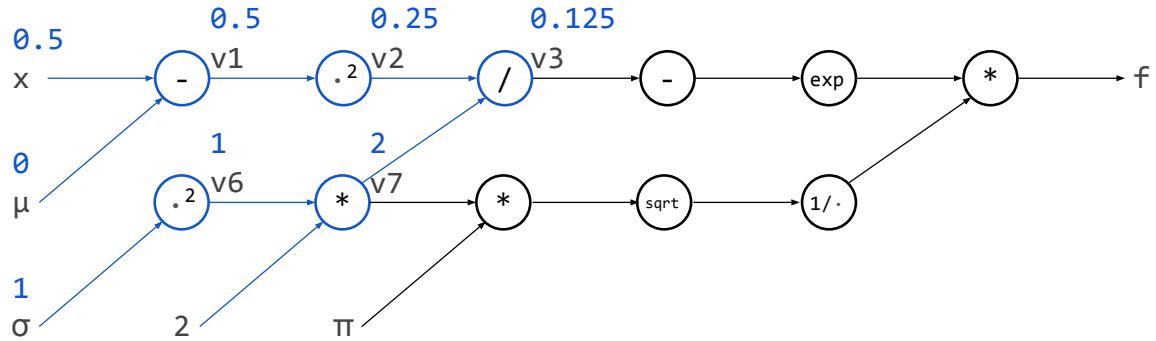
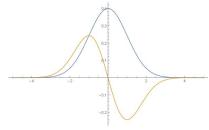
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



12

Normal PDF

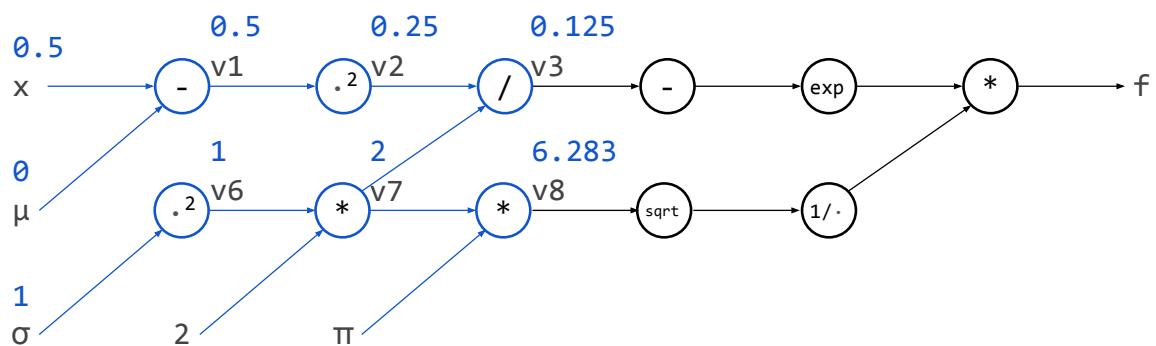
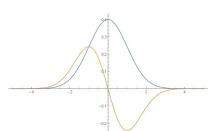
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



13

Normal PDF

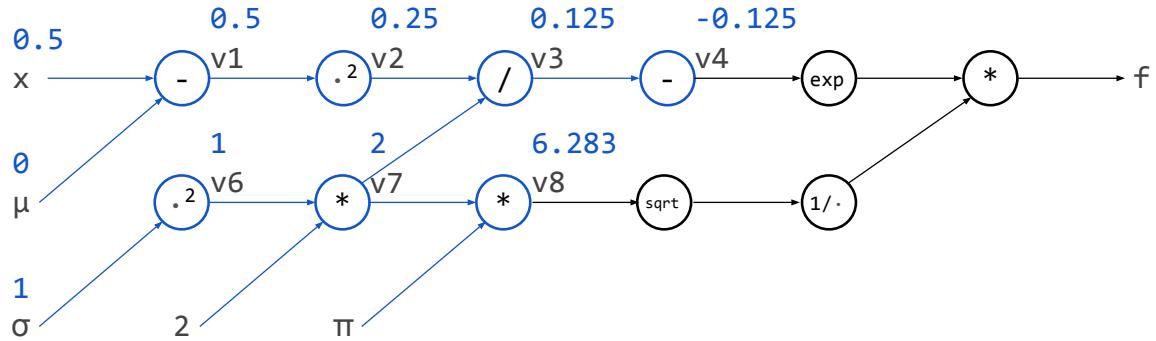
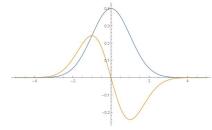
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



14

Normal PDF

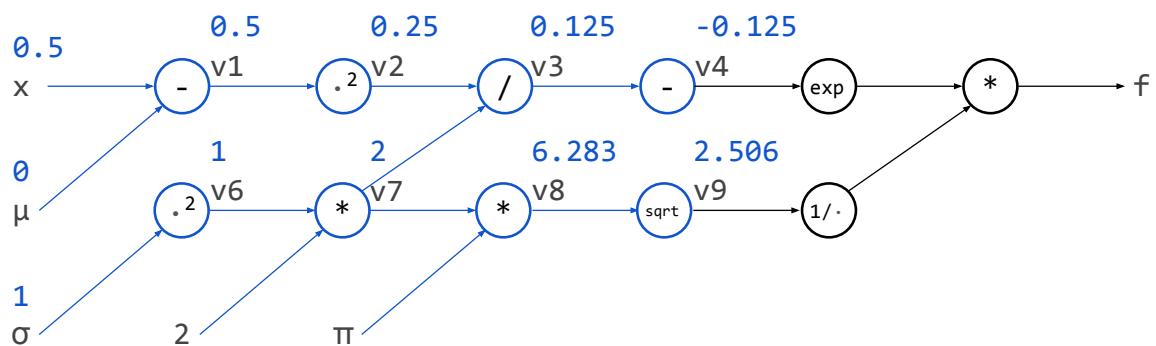
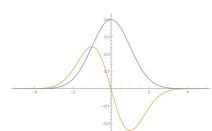
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



15

Normal PDF

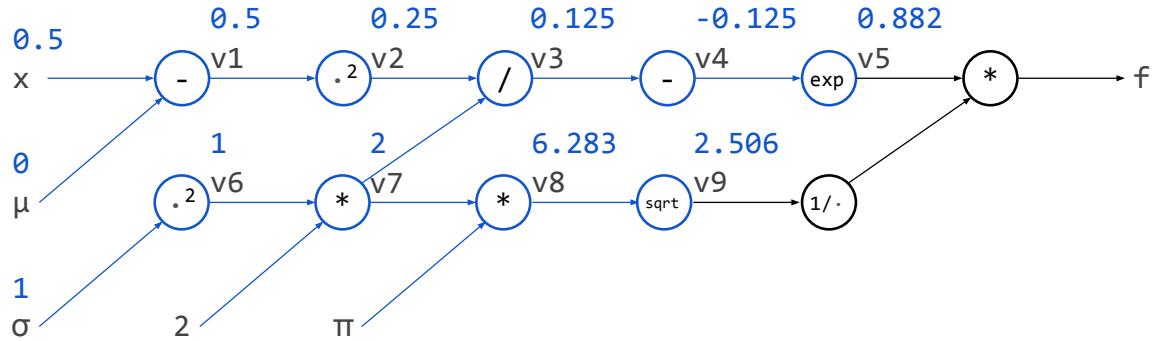
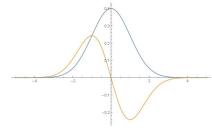
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



16

Normal PDF

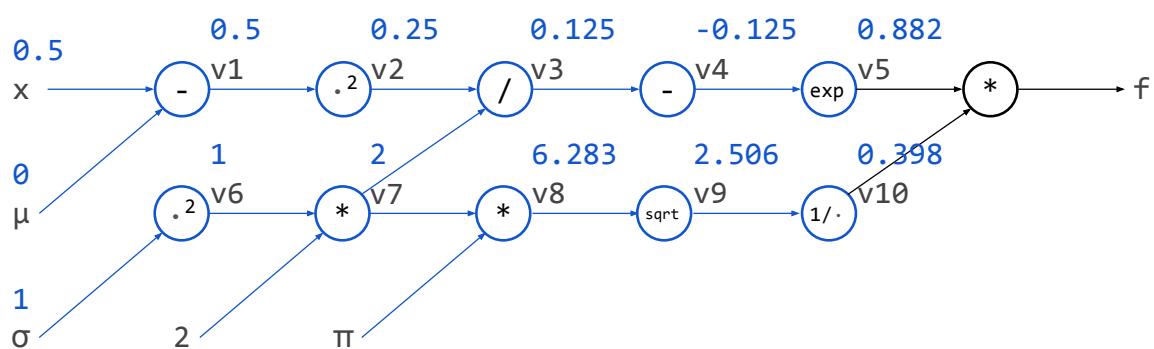
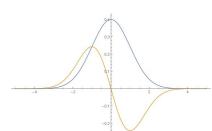
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



17

Normal PDF

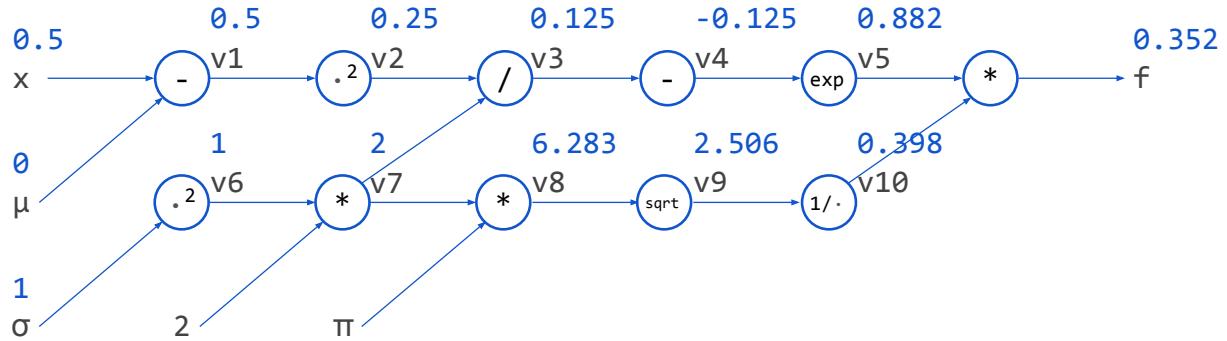
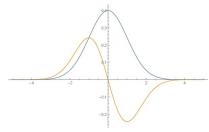
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



18

Normal PDF

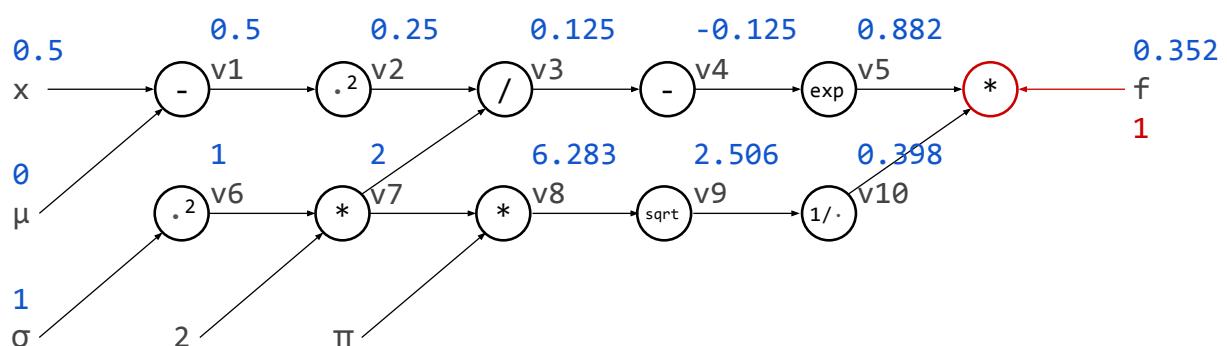
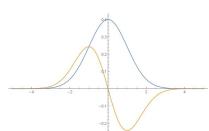
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



19

Normal PDF

$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

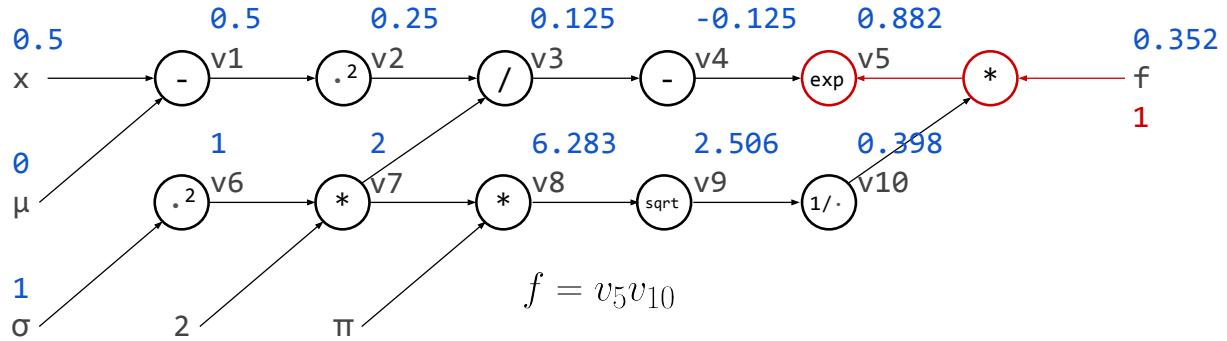
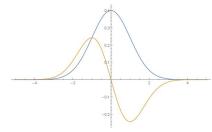


$$\frac{\partial f}{\partial f} = 1$$

20

Normal PDF

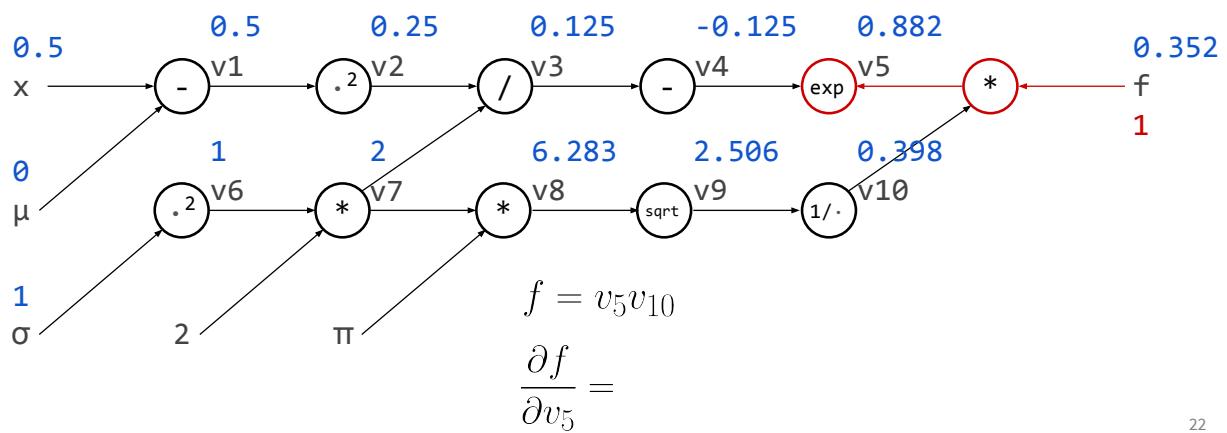
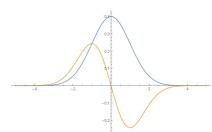
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



21

Normal PDF

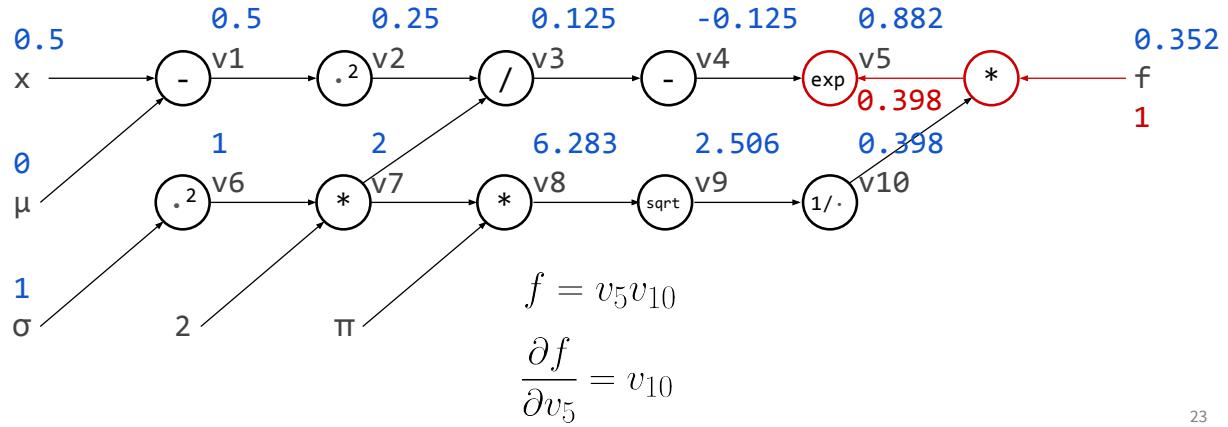
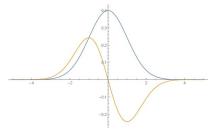
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



22

Normal PDF

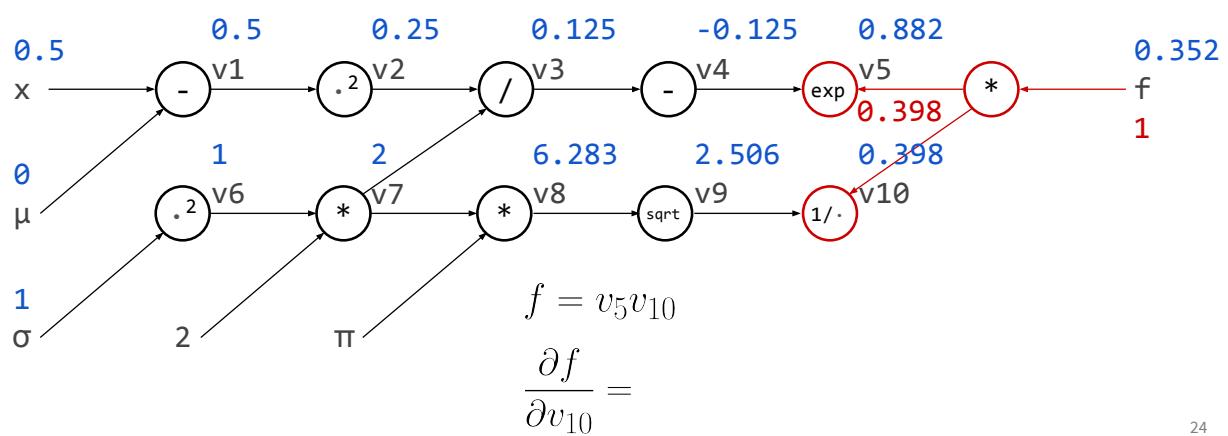
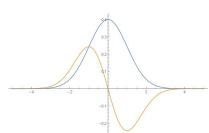
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



23

Normal PDF

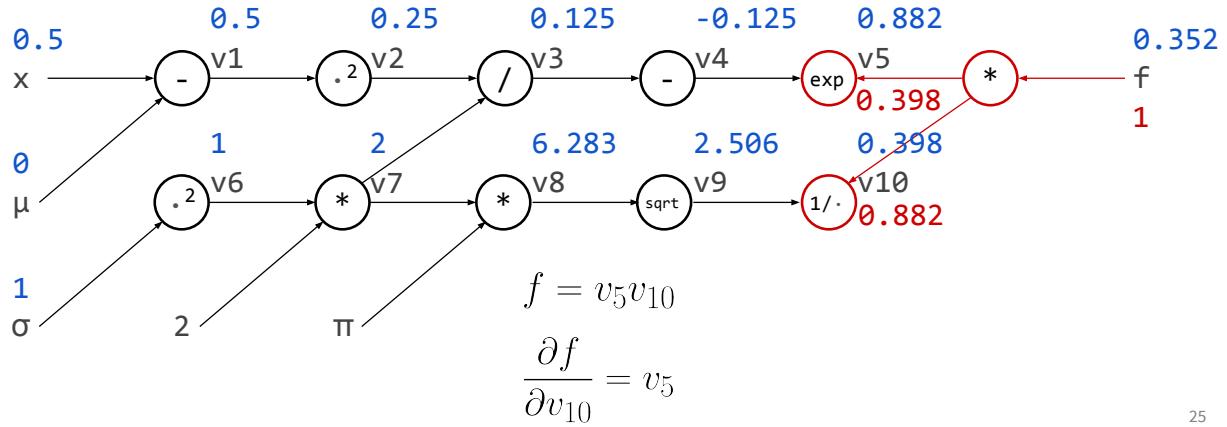
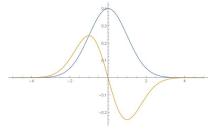
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



24

Normal PDF

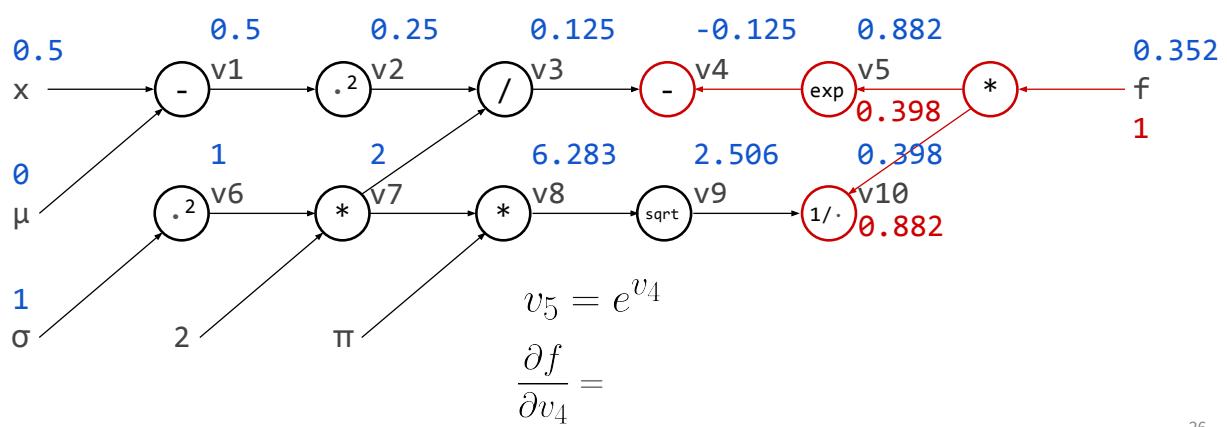
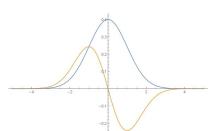
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



25

Normal PDF

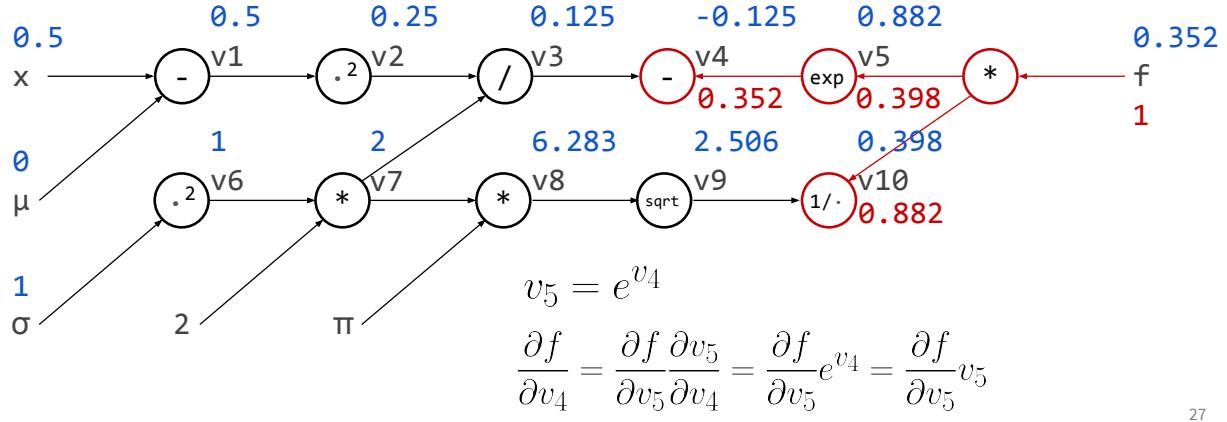
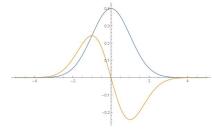
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



26

Normal PDF

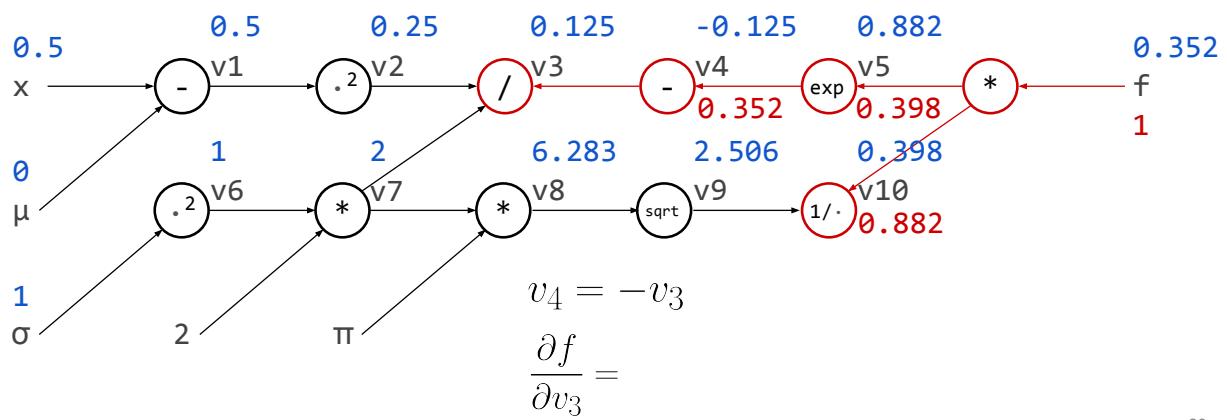
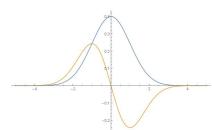
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



27

Normal PDF

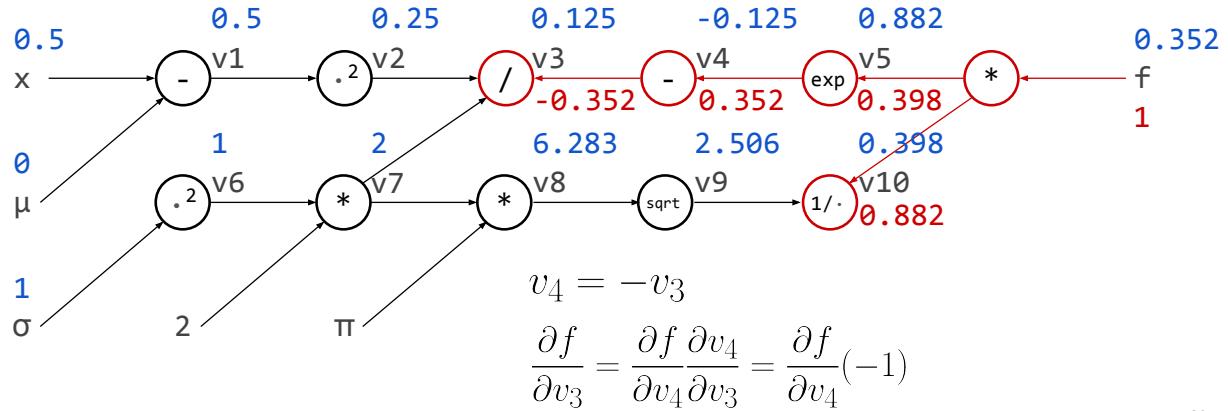
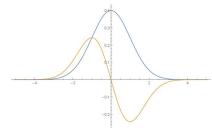
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



28

Normal PDF

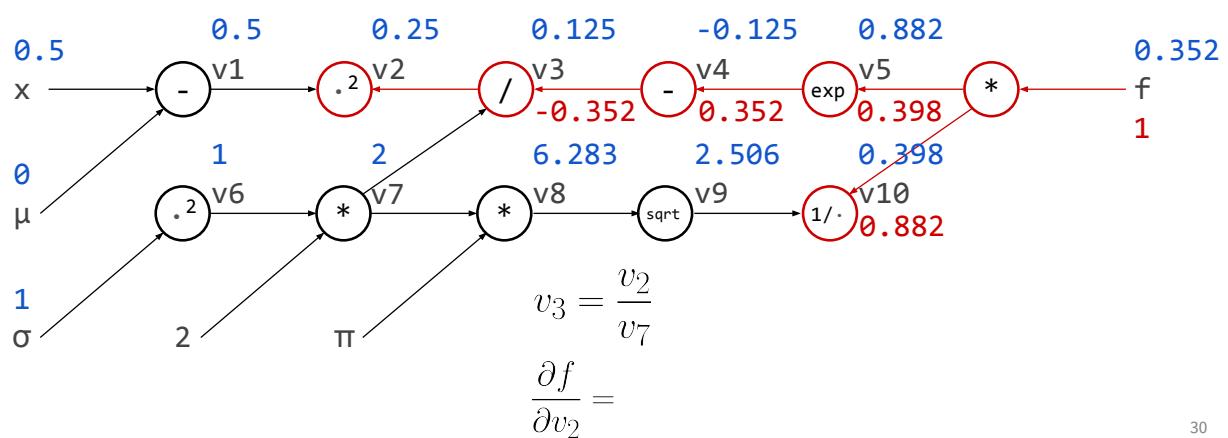
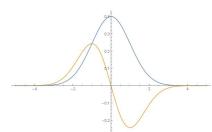
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



29

Normal PDF

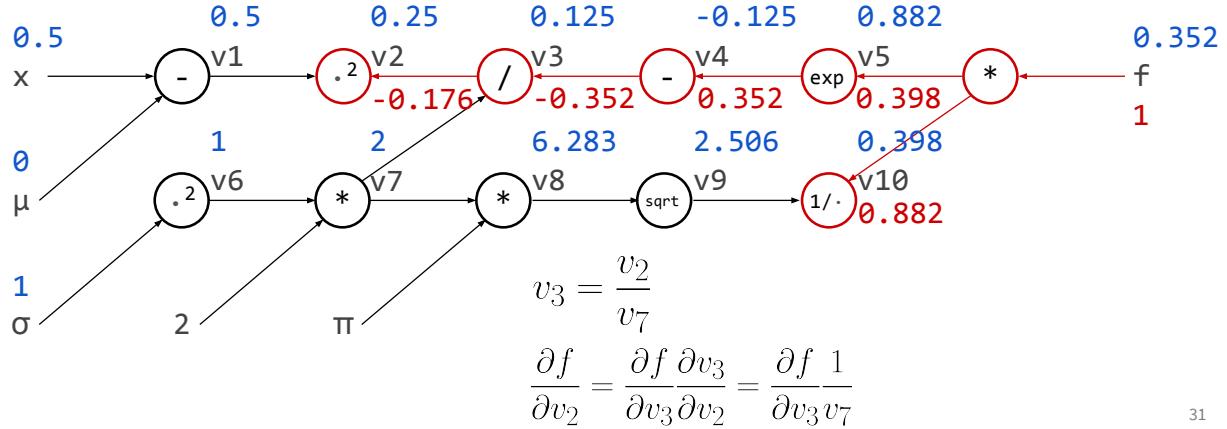
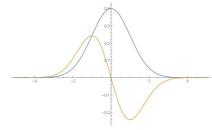
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



30

Normal PDF

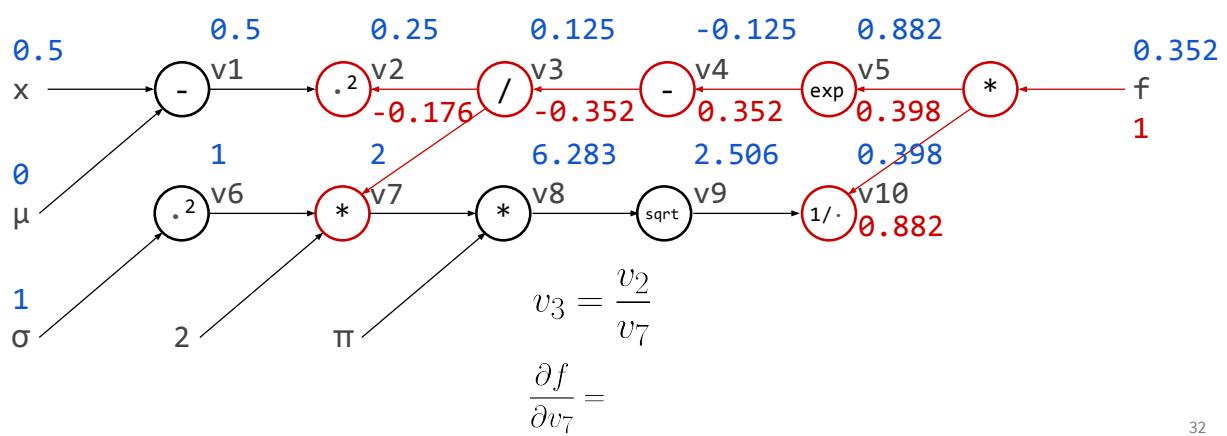
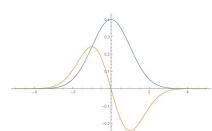
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



31

Normal PDF

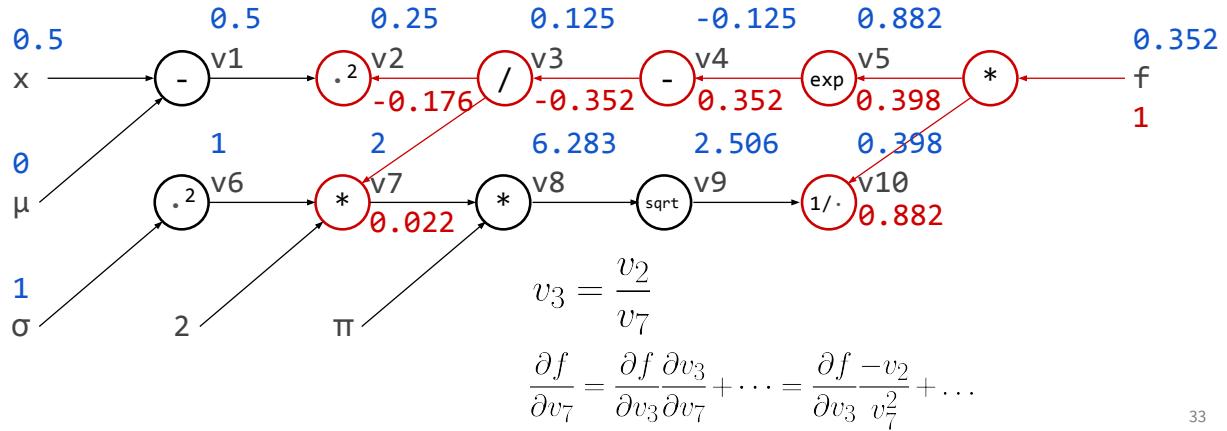
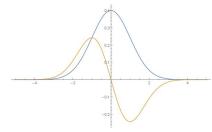
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



32

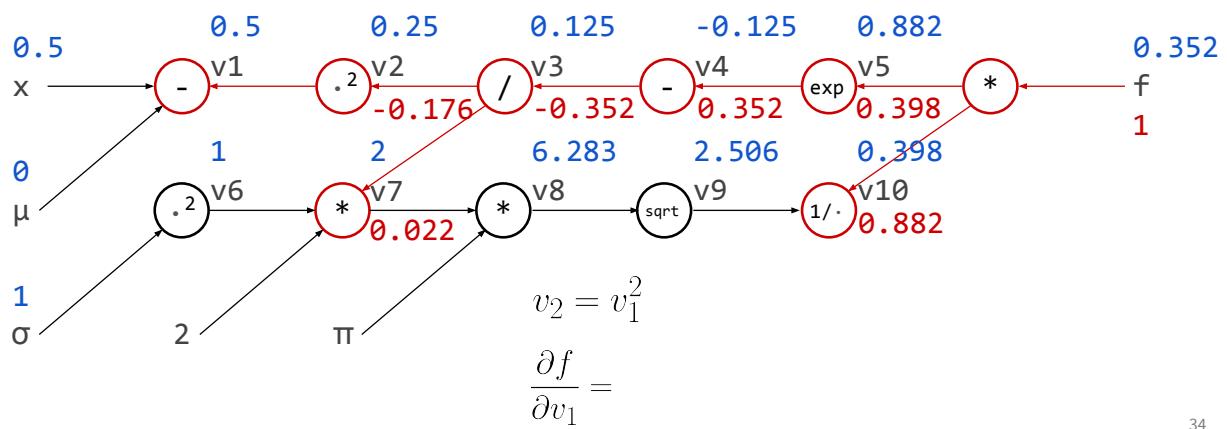
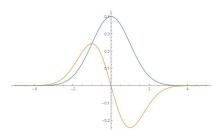
Normal PDF

$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



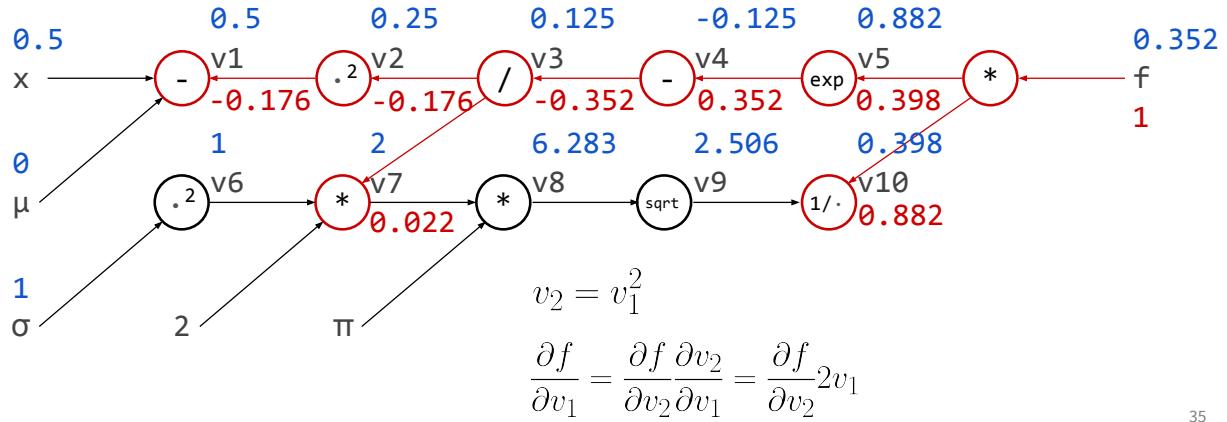
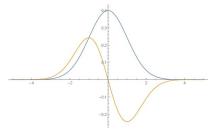
Normal PDF

$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



Normal PDF

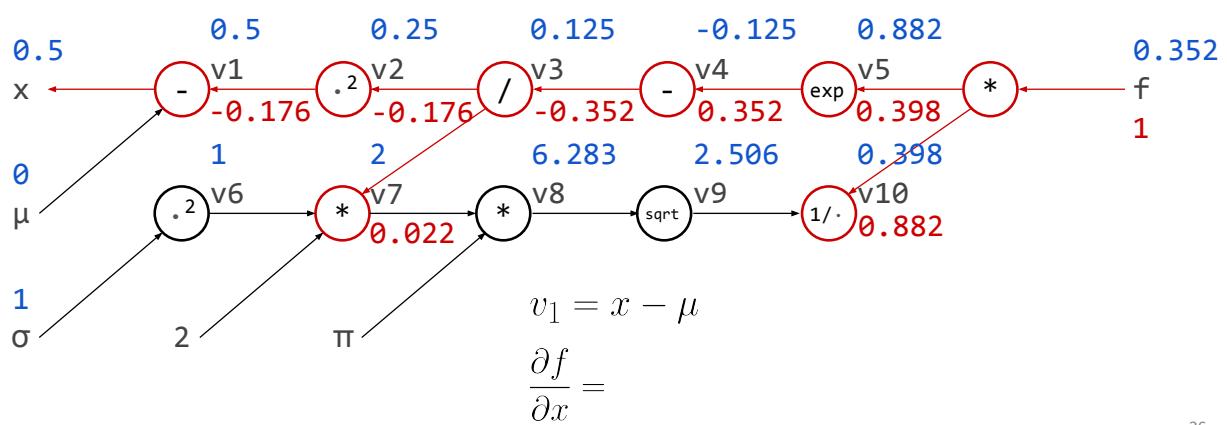
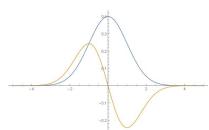
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



35

Normal PDF

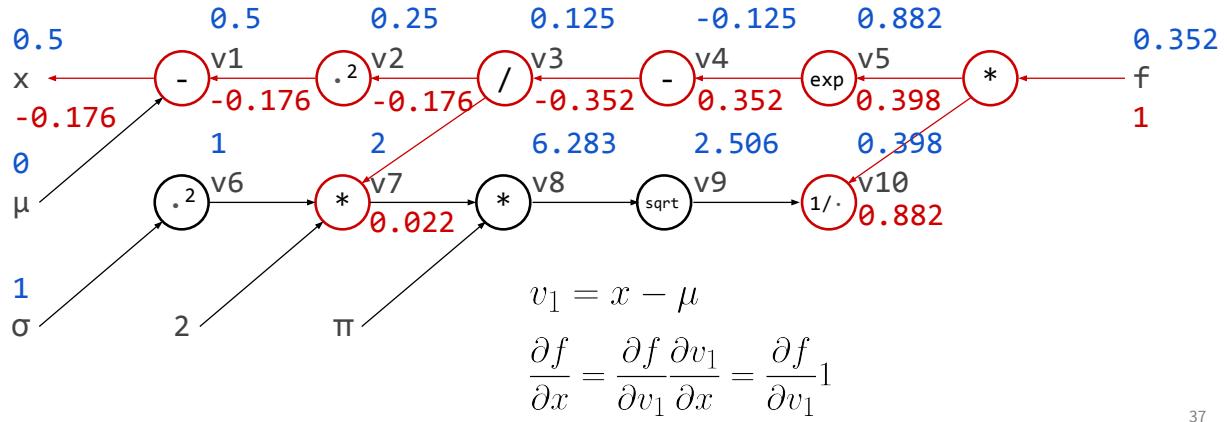
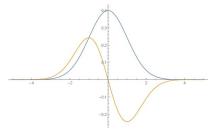
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



36

Normal PDF

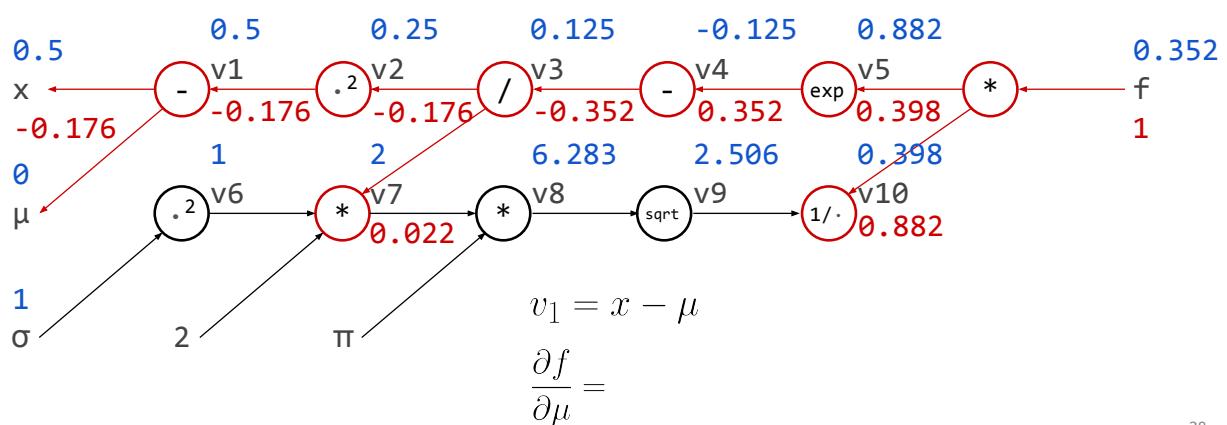
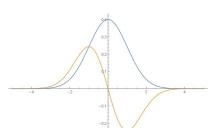
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



37

Normal PDF

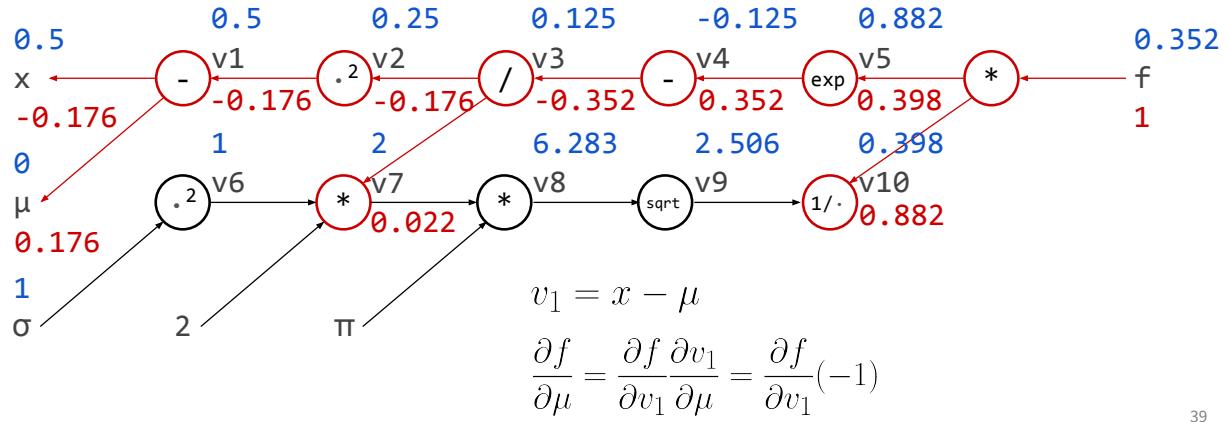
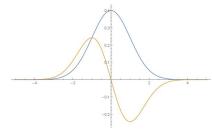
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



38

Normal PDF

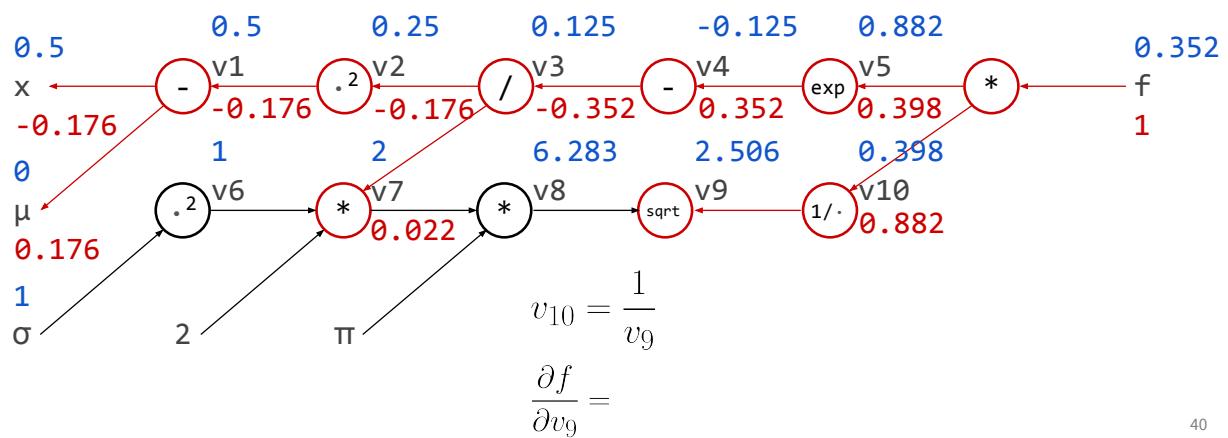
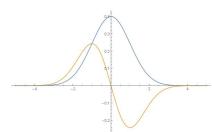
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



39

Normal PDF

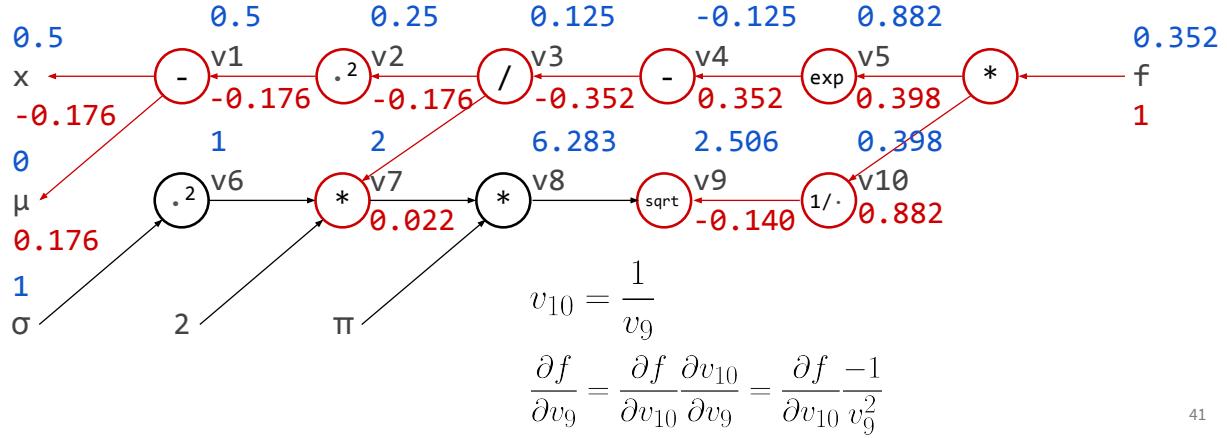
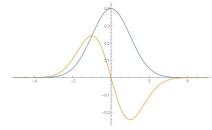
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



40

Normal PDF

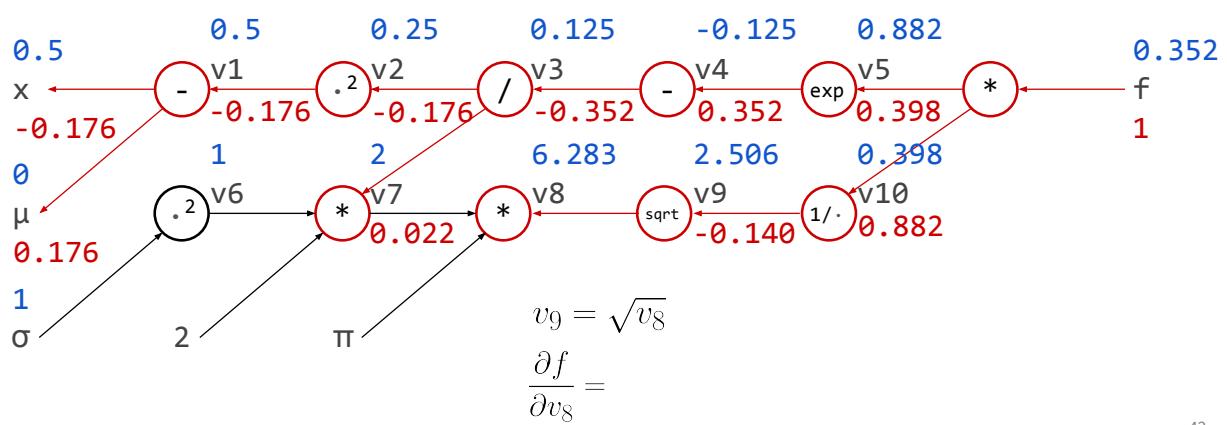
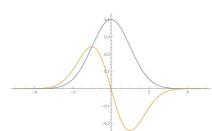
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



41

Normal PDF

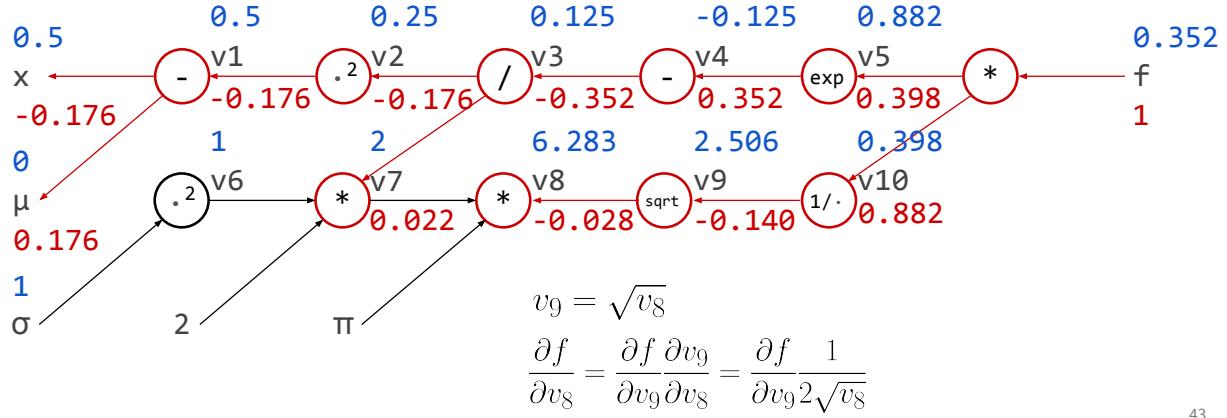
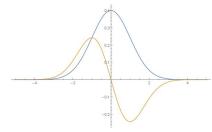
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



42

Normal PDF

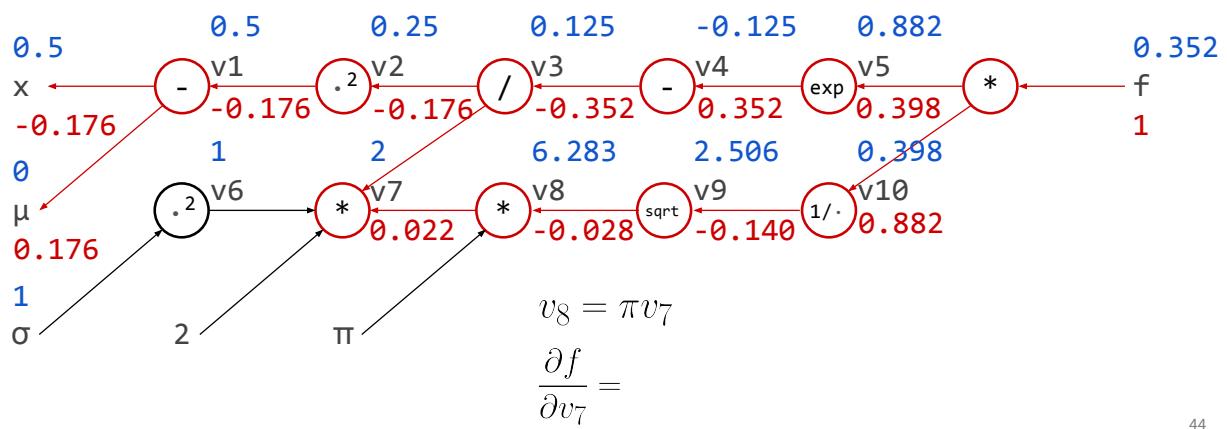
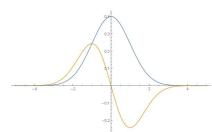
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



43

Normal PDF

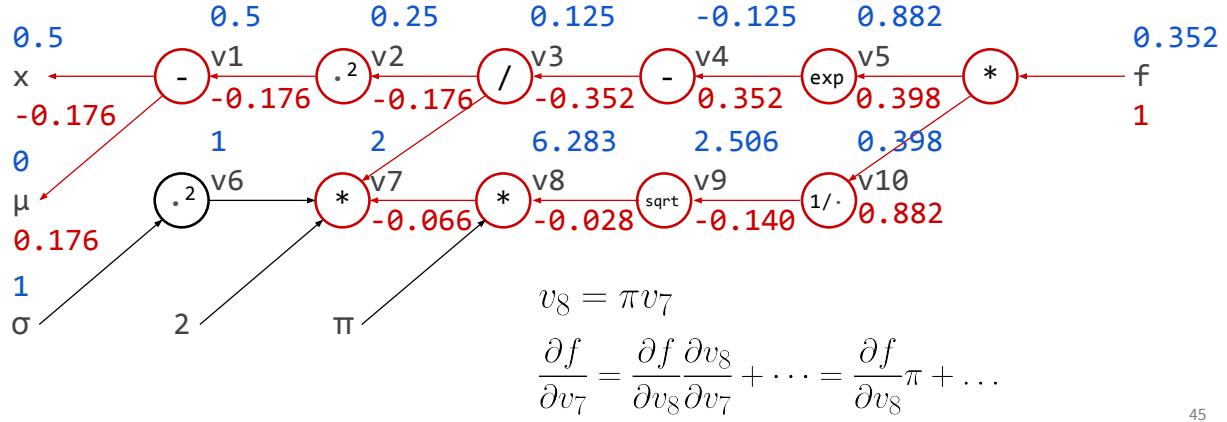
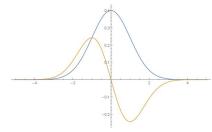
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



44

Normal PDF

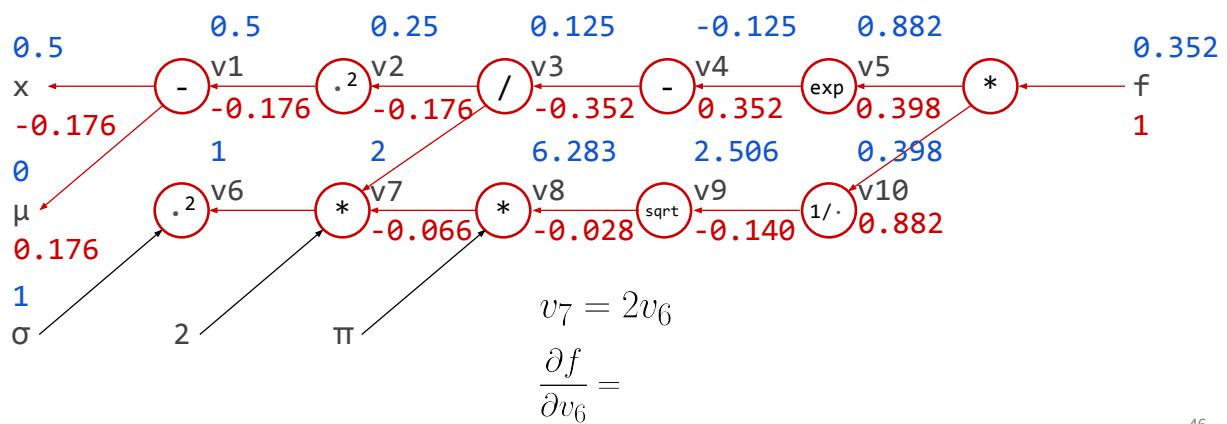
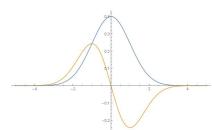
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



45

Normal PDF

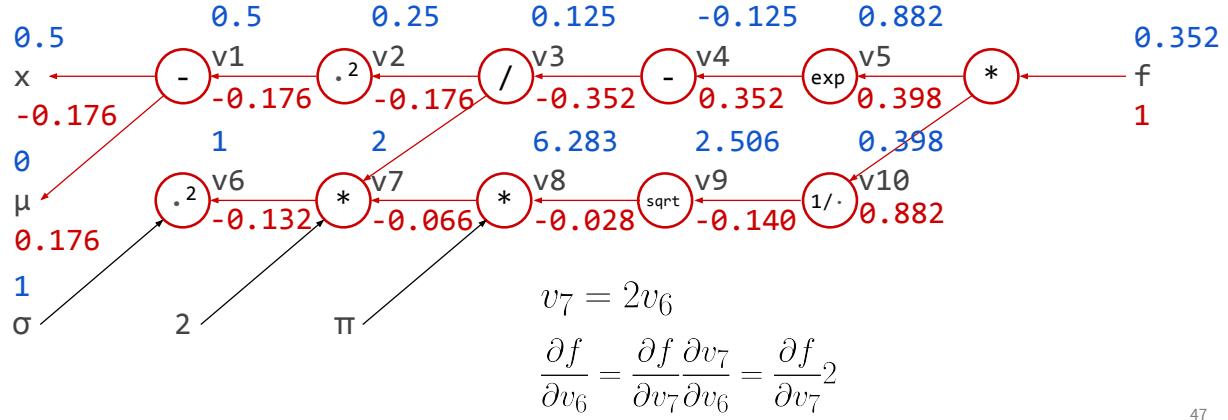
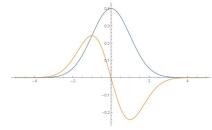
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



46

Normal PDF

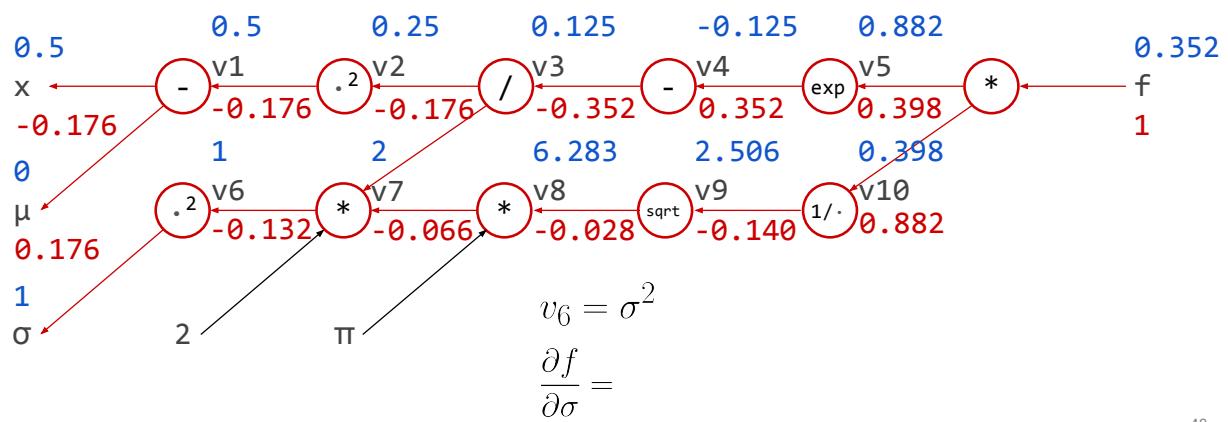
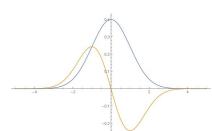
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



47

Normal PDF

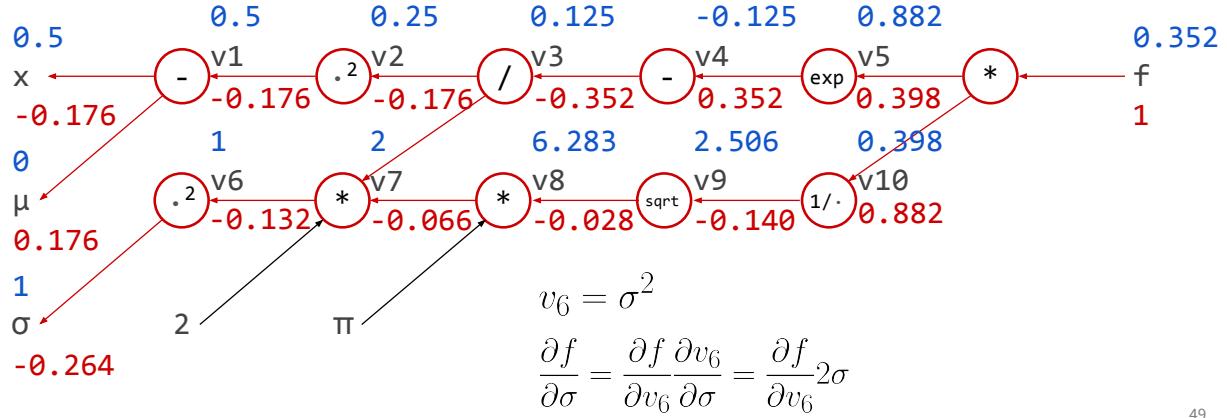
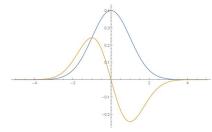
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



48

Normal PDF

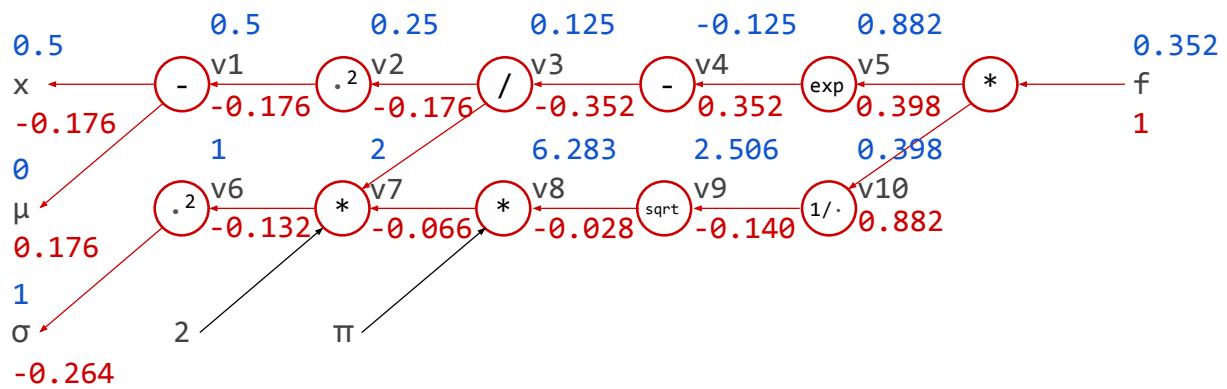
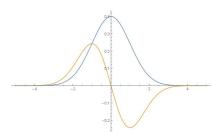
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



49

Normal PDF

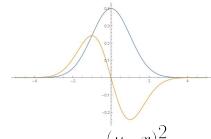
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



50

Normal PDF

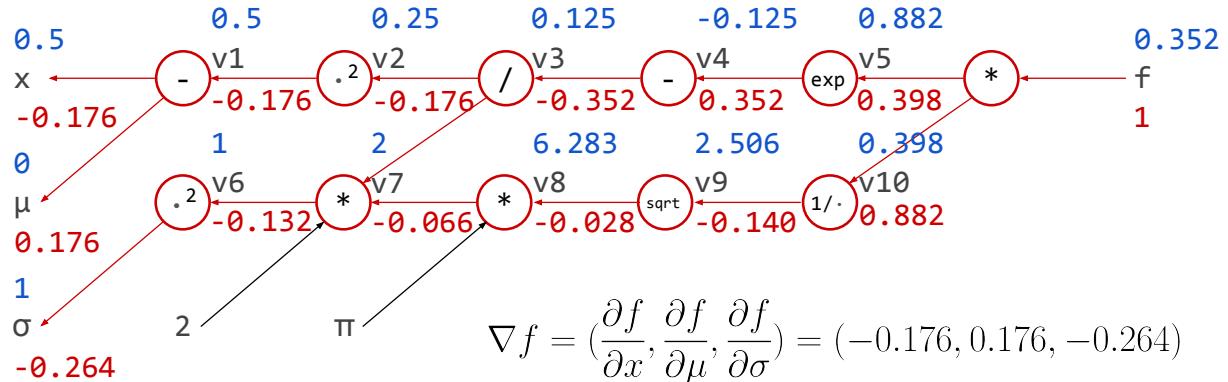
$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



$$\frac{\partial f}{\partial x} = \frac{(\mu - x)e^{-\frac{(\mu-x)^2}{2\sigma^2}}}{\sqrt{2\pi}\sigma^3}$$

$$\frac{\partial f}{\partial \mu} = \frac{(x - \mu)e^{-\frac{(x-\mu)^2}{2\sigma^2}}}{\sqrt{2\pi}\sigma^3}$$

$$\frac{\partial f}{\partial \sigma} = -\frac{(\sigma - x + \mu)(\sigma + x - \mu)e^{-\frac{(x-\mu)^2}{2\sigma^2}}}{\sqrt{2\pi}\sigma^4}$$



51

Implementation

52

Where does the graph come from?

Two main possibilities:

- **Static** computational graphs
Let the user define the graph as a data structure
- **Dynamic** computational graphs
Construct the graph automatically (general-purpose automatic differentiation)

53

Where does the graph come from?

Two main possibilities:

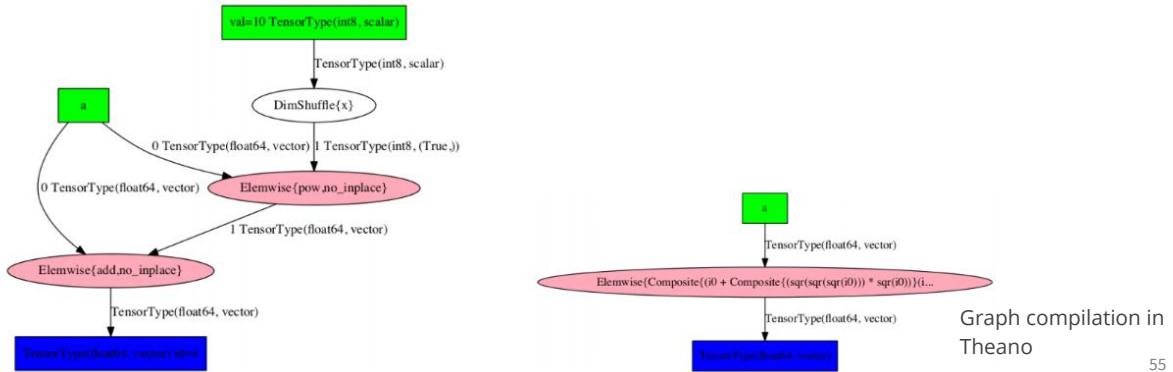
- **Static** computational graphs
Let the user define the graph as a data structure
“Define-and-run”
- **Dynamic** computational graphs
Construct the graph automatically (general-purpose automatic differentiation)
“Define-by-run”

54

Static graphs (define-and-run)

Prototypical examples: Theano, TensorFlow

- The user creates the graph using symbolic placeholders, with a mini-language
- Limited (and unintuitive) control flow and expressivity
- The graph gets “compiled” to take care of expression swell



Static graphs (define-and-run)

Prototypical examples: Theano, TensorFlow

Let's implement A^k

Static graphs (define-and-run)

Prototypical examples: Theano, TensorFlow

Let's implement A^k

Pure Python:

```
result = 1
for i in range(k):
    result = result * A
```

57

Static graphs (define-and-run)

Prototypical examples: Theano, TensorFlow

Let's implement A^k

Pure Python:

```
result = 1
for i in range(k):
    result = result * A
```

theano

```
import theano
import theano.tensor as T

k = T.iscalar("k")
A = T.vector("A")

# Symbolic description of a loop
result, updates = theano.scan(fn=lambda prior_result, A: prior_result * A,
                               outputs_info=T.ones_like(A),
                               non_sequences=A,
                               n_steps=k)
final_result = result[-1]

# Compiled function that returns A**k
power = theano.function(inputs=[A,k], outputs=final_result, updates=updates)
```

Dynamic graphs (define-by-run)

Prototypical examples: PyTorch (and TensorFlow eager execution)

- General-purpose autodiff, usually via operator overloading
- The user writes regular programs in host programming language
 - All language features (including control flow) are supported
- The graph is automatically constructed

59

Dynamic graphs (define-by-run)

Prototypical examples: PyTorch (and TensorFlow eager execution)

Let's implement A^k

60

Dynamic graphs (define-by-run)

Prototypical examples: PyTorch (and TensorFlow eager execution)

Let's implement A^k

Pure Python:

```
result = 1
for i in range(k):
    result = result * A
```

61

Dynamic graphs (define-by-run)

Prototypical examples: PyTorch (and TensorFlow eager execution)

Let's implement A^k

Pure Python:

```
result = 1
for i in range(k):
    result = result * A
```

 PYTORCH

```
import torch

result = torch.tensor(1)
for i in range(k):
    result = result * A

result.backward()
print(A.grad)
```

62

Where to implement

Many possibilities

- Interpreter-based
- Compiler-based
 - Source code transformation
 - Operator overloading

63

What to implement

Two main parts

Computational graph

- **Dynamically build the graph**
Side effect of forward evaluation or “non-standard interpretation”
- **Graph traversal algorithm**
The API to kickstart the backpropagation:
backward, grad, etc.

Derivatives

- **Rules of differentiation**
For all elementary numerical operations

Usually implemented on a custom numerical type, using operator overloading

64

Operator overloading on custom type

```
def eval_and_backprop(fun, x):
    y = fun(x)
    y.backprop(1.)
    return y, x._adjoint
```

```
class OpMul():
    def __init__(a, b):
        self._a = a
        self._b = b
        return a * b

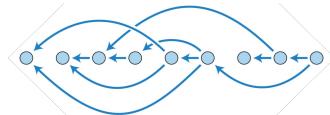
    def backprop(self, adjoint):
        self._a.backprop(adjoint * b)
        self._b.backprop(adjoint * a)
```

```
class Number():
    def __init__(op, fan_out=0):
        self._op = op
        self._adjoint = 0
        self._fan_out = fan_out

    def backprop(self, adjoint):
        self._adjoint += adjoint
        self._fan_out -= 1
        if self._fan_out == 0:
            self._op.backprop(self._adjoint)

    def __mul__(self, other):
        self._fan_out += 1
        other._fan_out += 1
        return Number(OpMul(self, other))
```

Graph with children pointing to parent(s)



65

Operator overloading on custom type

```
def eval_and_backprop(fun, x):
    y = fun(x)
    y.backprop(1.)
    return y, x._adjoint
```

```
class OpMul():
    def __init__(a, b):
        self._a = a
        self._b = b
        return a * b

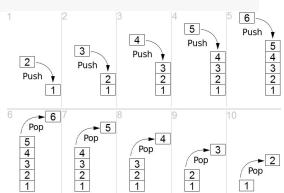
    def backprop(self, adjoint):
        self._a.backprop(adjoint * b)
        self._b.backprop(adjoint * a)
```

```
tape = []

class Number():
    def __init__(value):
        self._value = value
        self._adjoint = 0

    def backprop(self, adjoint):
        self._adjoint += adjoint

    def __mul__(self, other):
        global tape
        op = OpMul(self, other)
        tape.append(OpMul(self, other))
        return Number(self * other)
```



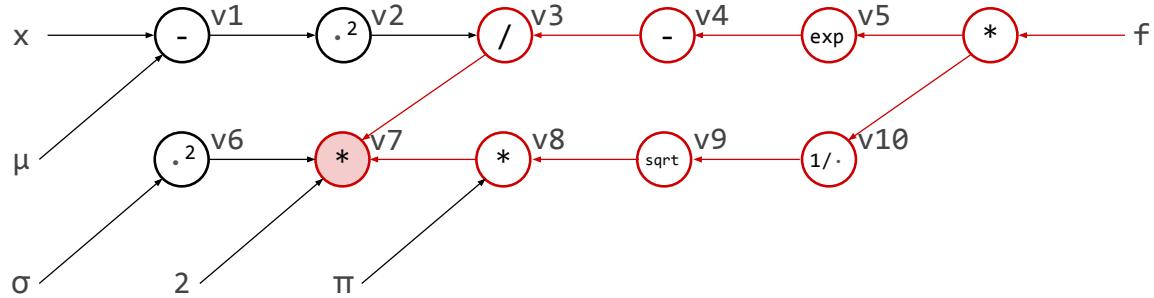
Global tape (stack)

- Forward: push in the order of evaluation
- Reverse: pop in the reverse order

66

Correctly handle fan-out

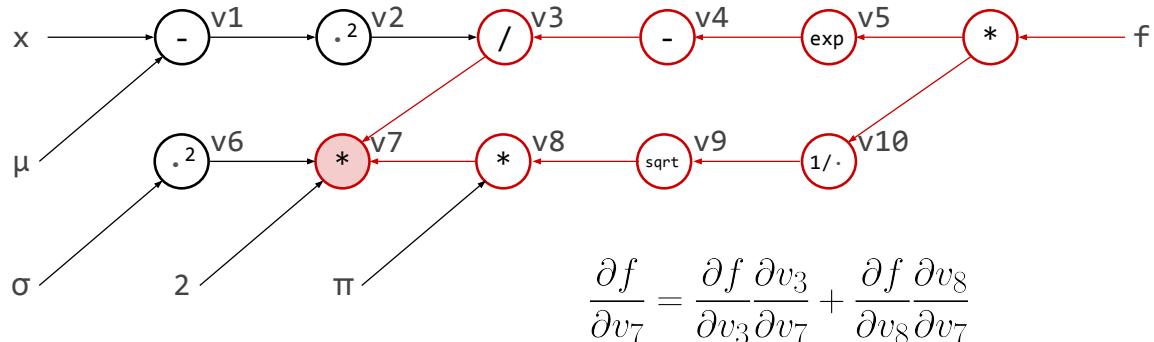
Fan-out: when a node is involved in multiple subsequent operations



67

Correctly handle fan-out

Fan-out: when a node is involved in multiple subsequent operations



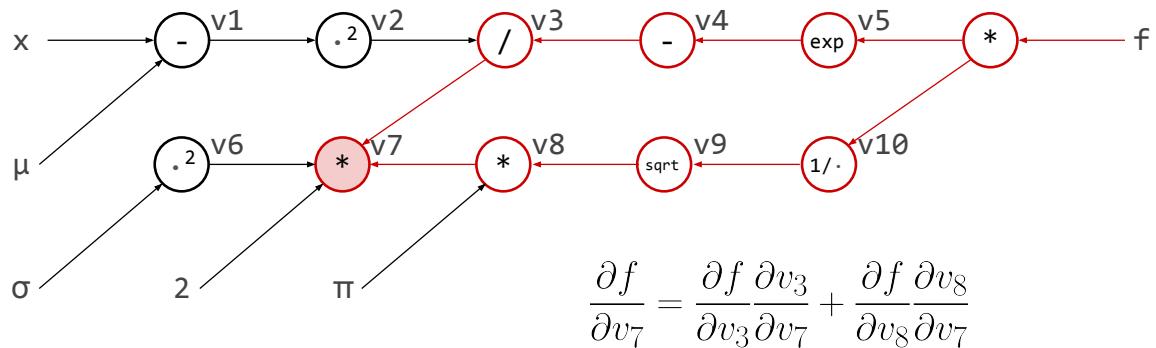
$$\frac{\partial f}{\partial v_7} = \frac{\partial f}{\partial v_3} \frac{\partial v_3}{\partial v_7} + \frac{\partial f}{\partial v_8} \frac{\partial v_8}{\partial v_7}$$

68

Correctly handle fan-out

Fan-out: when a node is involved in multiple subsequent operations

- Maintain a fan-out counter per node
- Don't propagate backward from a node until all derivatives coming to that node have arrived

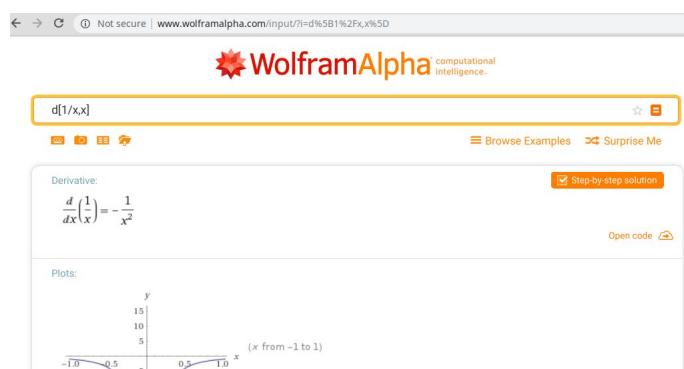


69

Check for correctness

Use numerical and symbolic differentiation to check individual rules and your chain rule implementation

$$\frac{\partial f(\mathbf{x})}{\partial x_i} = \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x} - h\mathbf{e}_i)}{2h} + O(h^2)$$



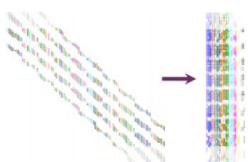
70

Advanced concepts

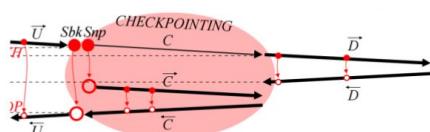
71

Some advanced concepts

- Nesting (e.g., forward-on-reverse), Hessian–vector products (Pearlmutter, 1994)
- Tape reduction and elimination (Naumann, 2004)
- Context-aware source-to-source transformation (Utke, 2004)
- Sparsity-aware autodiff by matrix coloring (Gebremedhin et al., 2013)



- Reverse mode checkpointing
(Dauvergne & Hascoet, 2006)



<https://github.com/openai/gradient-checkpointing>

Gruslys, A., Munos, R., Danihelka, I., Lanctot, M. and Graves, A.,
2016. Memory-efficient backpropagation through time. NIPS 2016⁷²

Summary

73

Summary

- Derivatives in machine learning
- How do we compute derivatives: manual, symbolic, numerical
- Automatic differentiation
 - Forward and reverse: $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ use reverse when $n \gg m$
 - Keep forward in mind if you need something more than backprop
- Computational graphs and propagation
- Implementation
 - Where does the graph come from?
 - Strategies and performance tips
- Advanced concepts
 - Nesting, higher-order derivatives
 - Checkpointing

74

References

- Baydin, A.G., Pearlmutter, B.A., Radul, A.A. and Siskind, J.M., 2017. Automatic differentiation in machine learning: a survey. *Journal of Machine Learning Research (JMLR)*, 18(153), pp.1-153.
- Baydin, Atilim Güneş, Barak A. Pearlmutter, and Jeffrey Mark Siskind. 2016. "Tricks from Deep Learning." In 7th International Conference on Algorithmic Differentiation, Christ Church Oxford, UK, September 12–15, 2016.
- Baydin, Atilim Güneş, Barak A. Pearlmutter, and Jeffrey Mark Siskind. 2016. "DiffSharp: An AD Library for .NET Languages." In 7th International Conference on Algorithmic Differentiation, Christ Church Oxford, UK, September 12–15, 2016.
- Baydin, Atilim Güneş, Robert Cornish, David Martínez Rubio, Mark Schmidt, and Frank Wood. 2018. "Online Learning Rate Adaptation with Hypergradient Descent." In Sixth International Conference on Learning Representations (ICLR), Vancouver, Canada, April 30 – May 3, 2018.
- Griewank, A. and Walther, A., 2008. *Evaluating derivatives: principles and techniques of algorithmic differentiation* (Vol. 105). SIAM.
- Nocedal, J. and Wright, S.J., 1999. *Numerical Optimization*. Springer.

75

Extra slides

76

Forward vs reverse summary

In the extreme $\mathbf{f} : \mathbb{R} \rightarrow \mathbb{R}^m$
use forward mode to evaluate

$$\left(\frac{\partial f_1}{\partial x}, \dots, \frac{\partial f_m}{\partial x} \right)$$

In the extreme $f : \mathbb{R}^n \rightarrow \mathbb{R}$
use reverse mode to evaluate

$$\nabla f(\mathbf{x}) = \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right)$$

In general $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ the Jacobian $\mathbf{J}_f(\mathbf{x}) \in \mathbb{R}^{m \times n}$ can be evaluated in
- $O(n \text{ time}(\mathbf{f}))$ with forward mode
- $O(m \text{ time}(\mathbf{f}))$ with reverse mode

Reverse performs better when $n \gg m$

77

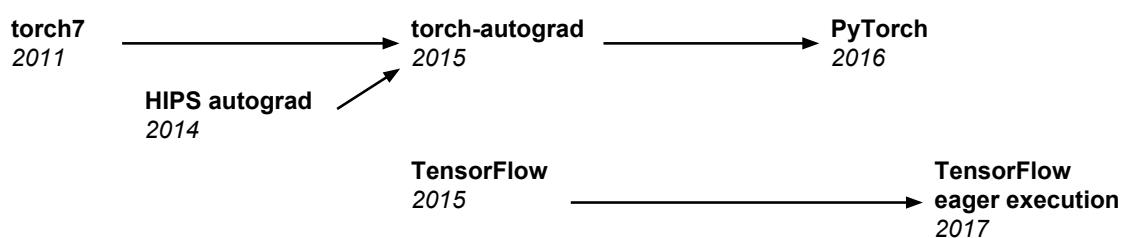
Current Landscape



Currently in progress: frameworks are in transition from
coarse-grained (module level) backprop

towards

fine-grained, general-purpose automatic differentiation

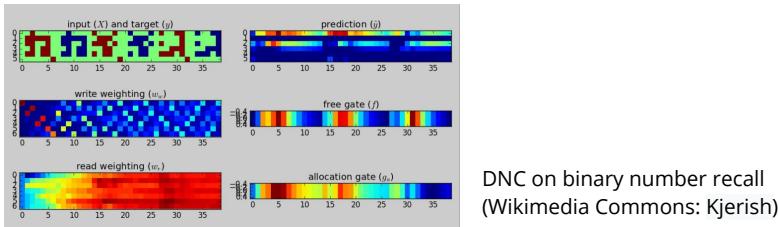


78

Current Landscape

A new mindset and workflow, enabling differentiable algorithmic elements

- Neural Turing Machine, Differentiable Neural Computer (Graves et al. 2014, 2016)
 - Can infer algorithms: copy, sort, recall
- Stack-augmented RNN (Joulin & Mikolov, 2015)
- End-to-end memory network (Sukhbaatar et al., 2015)
- Stack, queue, deque (Grefenstette et al., 2015)
- Discrete interfaces (Zaremba & Sutskever, 2015)



DNC on binary number recall
(Wikimedia Commons: Kjerish)

79

Current Landscape

General-purpose AD enables new libraries such as Pyro

Example:

Pyro supports stochastic recursion, higher-order functions, random control flow and runs stochastic variational inference enabled by PyTorch autograd infrastructure

PyTorch 2016 → Pyro 2017



```
In [7]: def geometric(p, t=None):
    if t is None:
        t = 0
    x = pyro.sample("x_{}".format(t), dist.bernoulli, p)
    if torch.equal(x.data, torch.zeros(1)):
        return x
    else:
        return x + geometric(p, t+1)

print(geometric(Variable(torch.Tensor([0.5]))))
```

Variable containing:
0
[torch.FloatTensor of size 1]

Dynamically generating random variables

80

DiffSharp

Op.	Value	Type signature	AD	Num.	Sym.
$f : \mathbb{R} \rightarrow \mathbb{R}$					
diff	f'	$(\mathbb{R} \rightarrow \mathbb{R}) \rightarrow \mathbb{R} \rightarrow \mathbb{R}$	X, F	A	X
diff'	(f, f')	$(\mathbb{R} \rightarrow \mathbb{R}) \rightarrow \mathbb{R} \rightarrow (\mathbb{R} \times \mathbb{R})$	X, F	A	X
diff2	f''	$(\mathbb{R} \rightarrow \mathbb{R}) \rightarrow \mathbb{R} \rightarrow \mathbb{R}$	X, F	A	X
diff2'	(f, f'')	$(\mathbb{R} \rightarrow \mathbb{R}) \rightarrow \mathbb{R} \rightarrow (\mathbb{R} \times \mathbb{R})$	X, F	A	X
diff2'',	(f, f', f'')	$(\mathbb{R} \rightarrow \mathbb{R}) \rightarrow \mathbb{R} \rightarrow (\mathbb{R} \times \mathbb{R} \times \mathbb{R})$	X, F	A	X
diffn	$f^{(n)}$	$\mathbb{N} \rightarrow (\mathbb{R} \rightarrow \mathbb{R}) \rightarrow \mathbb{R} \rightarrow \mathbb{R}$	X, F		X
diffn'	$(f, f^{(n)})$	$\mathbb{N} \rightarrow (\mathbb{R} \rightarrow \mathbb{R}) \rightarrow \mathbb{R} \rightarrow (\mathbb{R} \times \mathbb{R})$	X, F		X
$f : \mathbb{R}^n \rightarrow \mathbb{R}$					
grad	∇f	$(\mathbb{R}^n \rightarrow \mathbb{R}) \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^n$	X, R	A	X
grad'	$(f, \nabla f)$	$(\mathbb{R}^n \rightarrow \mathbb{R}) \rightarrow \mathbb{R}^n \rightarrow (\mathbb{R} \times \mathbb{R}^n)$	X, R	A	X
gradv	$\nabla f \cdot \mathbf{v}$	$(\mathbb{R}^n \rightarrow \mathbb{R}) \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}$	X, F	A	
gradv'	$(f, \nabla f \cdot \mathbf{v})$	$(\mathbb{R}^n \rightarrow \mathbb{R}) \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^n \rightarrow (\mathbb{R} \times \mathbb{R})$	X, F	A	
hessian	\mathbf{H}_f	$(\mathbb{R}^n \rightarrow \mathbb{R}) \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^{n \times n}$	X, R-F	A	X
hessian'	(f, \mathbf{H}_f)	$(\mathbb{R}^n \rightarrow \mathbb{R}) \rightarrow \mathbb{R}^n \rightarrow (\mathbb{R} \times \mathbb{R}^{n \times n})$	X, R-F	A	X
hessianv	$\mathbf{H}_f \mathbf{v}$	$(\mathbb{R}^n \rightarrow \mathbb{R}) \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^n$	X, F-R	A	
hessianv'	$(f, \mathbf{H}_f \mathbf{v})$	$(\mathbb{R}^n \rightarrow \mathbb{R}) \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^n \rightarrow (\mathbb{R} \times \mathbb{R}^n)$	X, F-R	A	
gradhessian	$(\nabla f, \mathbf{H}_f)$	$(\mathbb{R}^n \rightarrow \mathbb{R}) \rightarrow \mathbb{R}^n \rightarrow (\mathbb{R}^n \times \mathbb{R}^{n \times n})$	X, R-F	A	X
gradhessian'	$(f, \nabla f, \mathbf{H}_f)$	$(\mathbb{R}^n \rightarrow \mathbb{R}) \rightarrow \mathbb{R}^n \rightarrow (\mathbb{R} \times \mathbb{R}^n \times \mathbb{R}^{n \times n})$	X, R-F	A	X
gradhessianv	$(\nabla f \cdot \mathbf{v}, \mathbf{H}_f \mathbf{v})$	$(\mathbb{R}^n \rightarrow \mathbb{R}) \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^n \rightarrow (\mathbb{R} \times \mathbb{R}^n)$	X, F-R	A	
gradhessianv'	$(f, \nabla f \cdot \mathbf{v}, \mathbf{H}_f \mathbf{v})$	$(\mathbb{R}^n \rightarrow \mathbb{R}) \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^n \rightarrow (\mathbb{R} \times \mathbb{R} \times \mathbb{R}^n)$	X, F-R	A	
laplacian	$\text{tr}(\mathbf{H}_f)$	$(\mathbb{R}^n \rightarrow \mathbb{R}) \rightarrow \mathbb{R} \rightarrow \mathbb{R}$	X, R-F	A	X
laplacian'	$(f, \text{tr}(\mathbf{H}_f))$	$(\mathbb{R}^n \rightarrow \mathbb{R}) \rightarrow \mathbb{R}^n \rightarrow (\mathbb{R} \times \mathbb{R})$	X, R-F	A	X
$\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$					
jacobian	$\mathbf{J}_{\mathbf{f}}$	$(\mathbb{R}^n \rightarrow \mathbb{R}^m) \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^{m \times n}$	X, F/R	A	X
jacobian'	$(\mathbf{f}, \mathbf{J}_{\mathbf{f}})$	$(\mathbb{R}^n \rightarrow \mathbb{R}^m) \rightarrow \mathbb{R}^n \rightarrow (\mathbb{R}^m \times \mathbb{R}^{m \times n})$	X, F/R	A	X
jacobianv	$\mathbf{J}_{\mathbf{f}} \mathbf{v}$	$(\mathbb{R}^n \rightarrow \mathbb{R}^m) \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^m$	X, F	A	
jacobianv'	$(\mathbf{f}, \mathbf{J}_{\mathbf{f}} \mathbf{v})$	$(\mathbb{R}^n \rightarrow \mathbb{R}^m) \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^n \rightarrow (\mathbb{R}^m \times \mathbb{R}^m)$	X, F	A	
jacobianT	$\mathbf{J}_{\mathbf{f}}^T$	$(\mathbb{R}^n \rightarrow \mathbb{R}^m) \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^{n \times m}$	X, F/R	A	X
jacobianT'	$(\mathbf{f}, \mathbf{J}_{\mathbf{f}}^T)$	$(\mathbb{R}^n \rightarrow \mathbb{R}^m) \rightarrow \mathbb{R}^n \rightarrow (\mathbb{R}^m \times \mathbb{R}^{n \times m})$	X, F/R	A	X
jacobianTv	$\mathbf{J}_{\mathbf{f}}^T \mathbf{v}$	$(\mathbb{R}^n \rightarrow \mathbb{R}^m) \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^m \rightarrow \mathbb{R}^n$	X, R		
jacobianTv'	$(\mathbf{f}, \mathbf{J}_{\mathbf{f}}^T \mathbf{v})$	$(\mathbb{R}^n \rightarrow \mathbb{R}^m) \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^m \rightarrow (\mathbb{R}^m \times \mathbb{R}^n)$	X, R		
jacobianTv'',	$(\mathbf{f}, \mathbf{J}_{\mathbf{f}}^T(\cdot))$	$(\mathbb{R}^n \rightarrow \mathbb{R}^m) \rightarrow \mathbb{R}^n \rightarrow (\mathbb{R}^m \times (\mathbb{R}^m \rightarrow \mathbb{R}^n))$	X, R		
curl	$\nabla \times \mathbf{f}$	$(\mathbb{R}^3 \rightarrow \mathbb{R}^3) \rightarrow \mathbb{R}^3 \rightarrow \mathbb{R}^3$	X, F	A	X
curl'	$(\mathbf{f}, \nabla \times \mathbf{f})$	$(\mathbb{R}^3 \rightarrow \mathbb{R}^3) \rightarrow \mathbb{R}^3 \rightarrow (\mathbb{R}^3 \times \mathbb{R}^3)$	X, F	A	X
div	$\nabla \cdot \mathbf{f}$	$(\mathbb{R}^n \rightarrow \mathbb{R}^n) \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}$	X, F	A	X
div'	$(\mathbf{f}, \nabla \cdot \mathbf{f})$	$(\mathbb{R}^n \rightarrow \mathbb{R}^n) \rightarrow \mathbb{R}^n \rightarrow (\mathbb{R}^n \times \mathbb{R})$	X, F	A	X
curldiv	$(\nabla \times \mathbf{f}, \nabla \cdot \mathbf{f})$	$(\mathbb{R}^3 \rightarrow \mathbb{R}^3) \rightarrow \mathbb{R}^3 \rightarrow (\mathbb{R}^3 \times \mathbb{R})$	X, F	A	X
curldiv'	$(\mathbf{f}, \nabla \times \mathbf{f}, \nabla \cdot \mathbf{f})$	$(\mathbb{R}^3 \rightarrow \mathbb{R}^3) \rightarrow \mathbb{R}^3 \rightarrow (\mathbb{R}^3 \times \mathbb{R}^3 \times \mathbb{R})$	X, F	A	X