

# Advanced Machine Learning

## Lecture 7: Deep neural networks

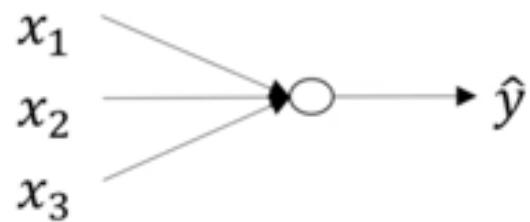
Sandjai Bhulai  
Vrije Universiteit Amsterdam

s.bhulai@vu.nl  
26 September 2023

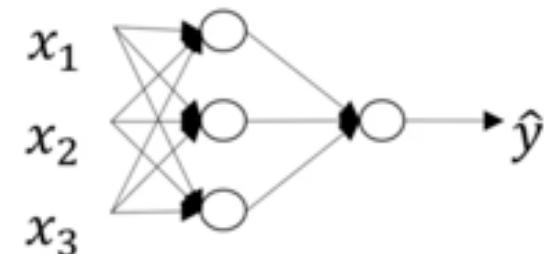
# Deep neural networks

Advanced Machine Learning

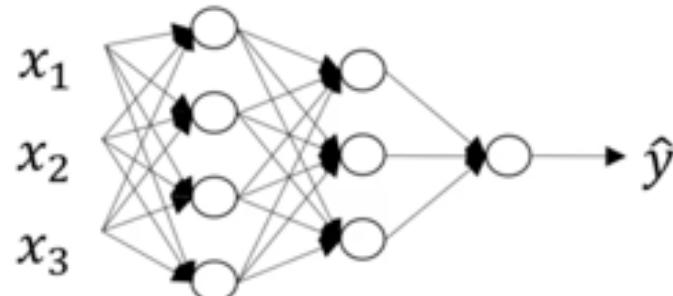
# Deep neural networks



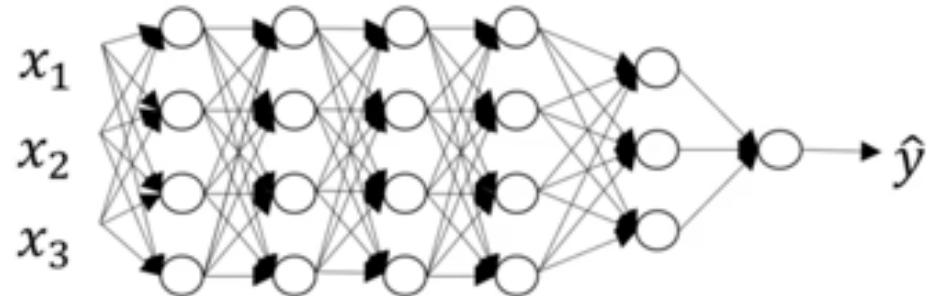
logistic regression



1 hidden layer

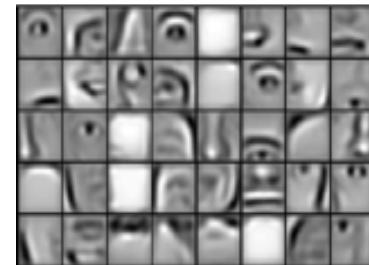
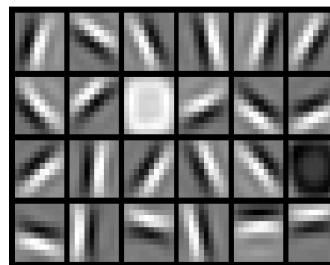
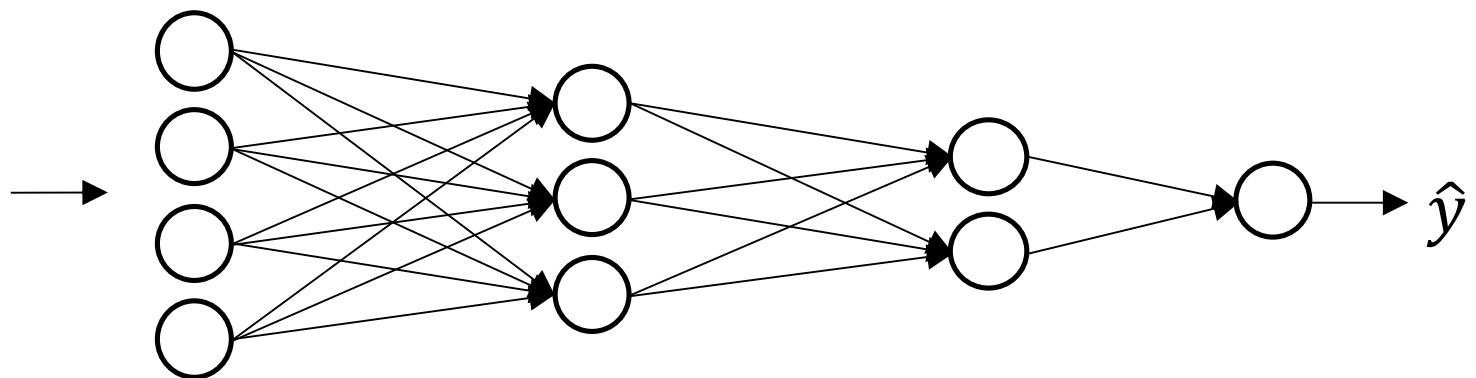


2 hidden layers



5 hidden layers

# Deep neural networks



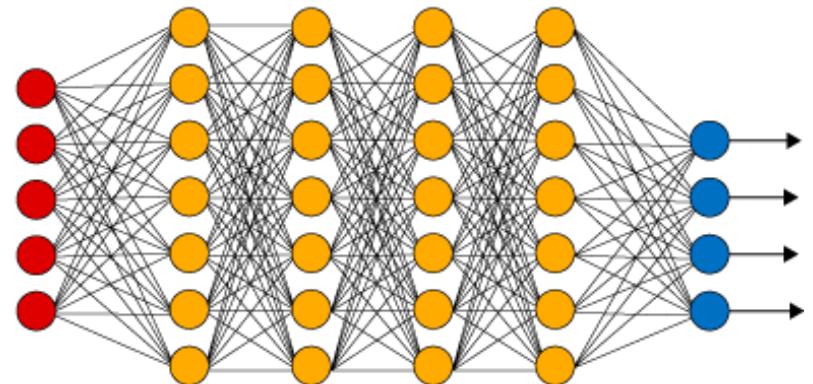
# Deep neural networks

- Why deep neural networks?
  - > Informally: There are many functions you can compute with a “small” -layer deep neural network that shallower networks require exponentially more hidden units to compute
  - > Example:  $x_1 \text{ XOR } x_2 \text{ XOR } x_3 \text{ XOR } x_4 \dots$   
 $O(\log n)$  versus  $O(2^n)$

# Deep neural networks

Building blocks for forward propagation:

- Consider layer  $l$
- Input to layer  $l$  is  $a^{[l-1]}$
- Output  $a^{[l]}$  given by:



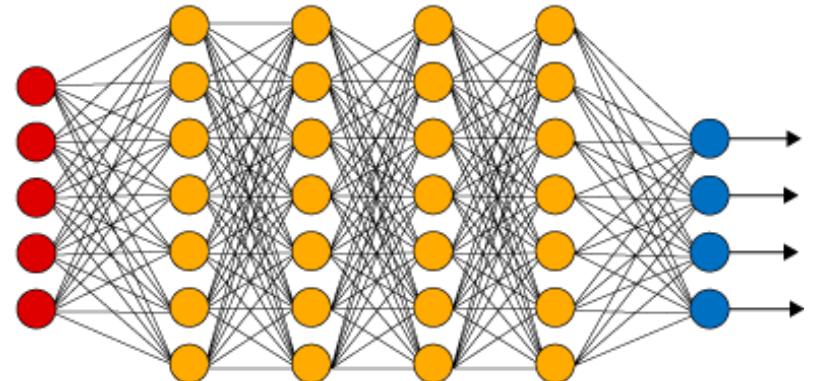
$$z^{[l]} = W^{[l]}a^{[l-1]} + b^{[l]}$$

$$a^{[l]} = g^{[l]}(z^{[l]})$$

# Deep neural networks

Building blocks for backward propagation:

- Consider layer  $l$
- Input to layer  $l$  is  $\mathbf{d}a^{[l]}$
- Output is given by:



$$\mathbf{d}z^{[l]} = \mathbf{d}a^{[l]} \cdot g^{[l]}'(z^{[l]})$$

$$\mathbf{d}W^{[l]} = \mathbf{d}z^{[l]} \cdot a^{[l-1]}$$

$$\mathbf{d}b^{[l]} = \mathbf{d}z^{[l]}$$

$$\mathbf{d}a^{[l-1]} = (W^{[l]})^\top \mathbf{d}z^{[l]}$$

# Deep neural networks

## Forward and backward propagation

```
Z[1] = W[1]X + b[1]
A[1] = g[1](Z[1])
Z[2] = W[2]A[1] + b[2]
A[2] = g[2](Z[2])
:
A[L] = g[L](Z[L]) = Ŷ
```

$$dZ^{[L]} = A^{[L]} - Y$$

$$dW^{[L]} = \frac{1}{m} dZ^{[L]} (A^{[L]})^\top$$

$$db^{[L]} = \frac{1}{m} \text{np.sum}(dZ^{[L]}, \text{axis} = 1, \text{keepdims} = \text{True})$$

$$dZ^{[L-1]} = d(W^{[L]})^\top dZ^{[L]} g'^{[L]}(Z^{[L-1]})$$

:

$$dZ^{[1]} = d(W^{[L]})^\top dZ^{[2]} g'^{[1]}(Z^{[1]})$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} (A^{[1]})^\top$$

$$db^{[1]} = \frac{1}{m} \text{np.sum}(dZ^{[1]}, \text{axis} = 1, \text{keepdims} = \text{True})$$

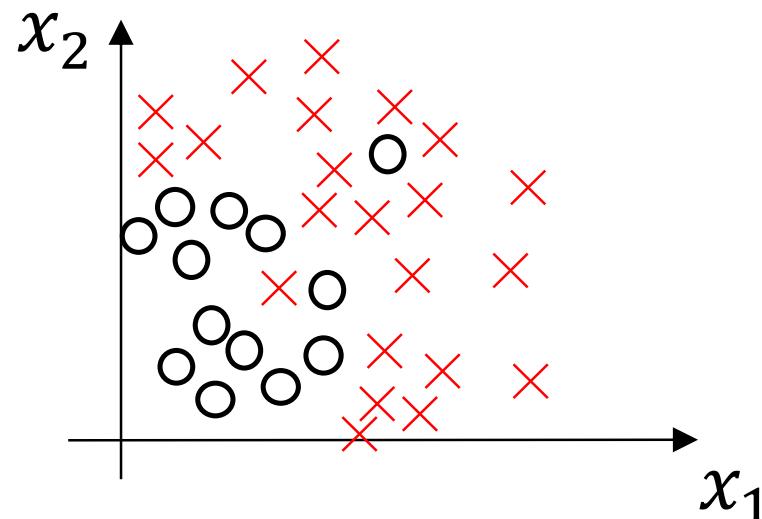
# Practical aspects of DNNs

- Deep neural networks have a lot of hyperparameters
  - > # layers
  - > # hidden units
  - > Learning rates
  - > Activation functions
  - > ...
- Need for: train, development, and test set

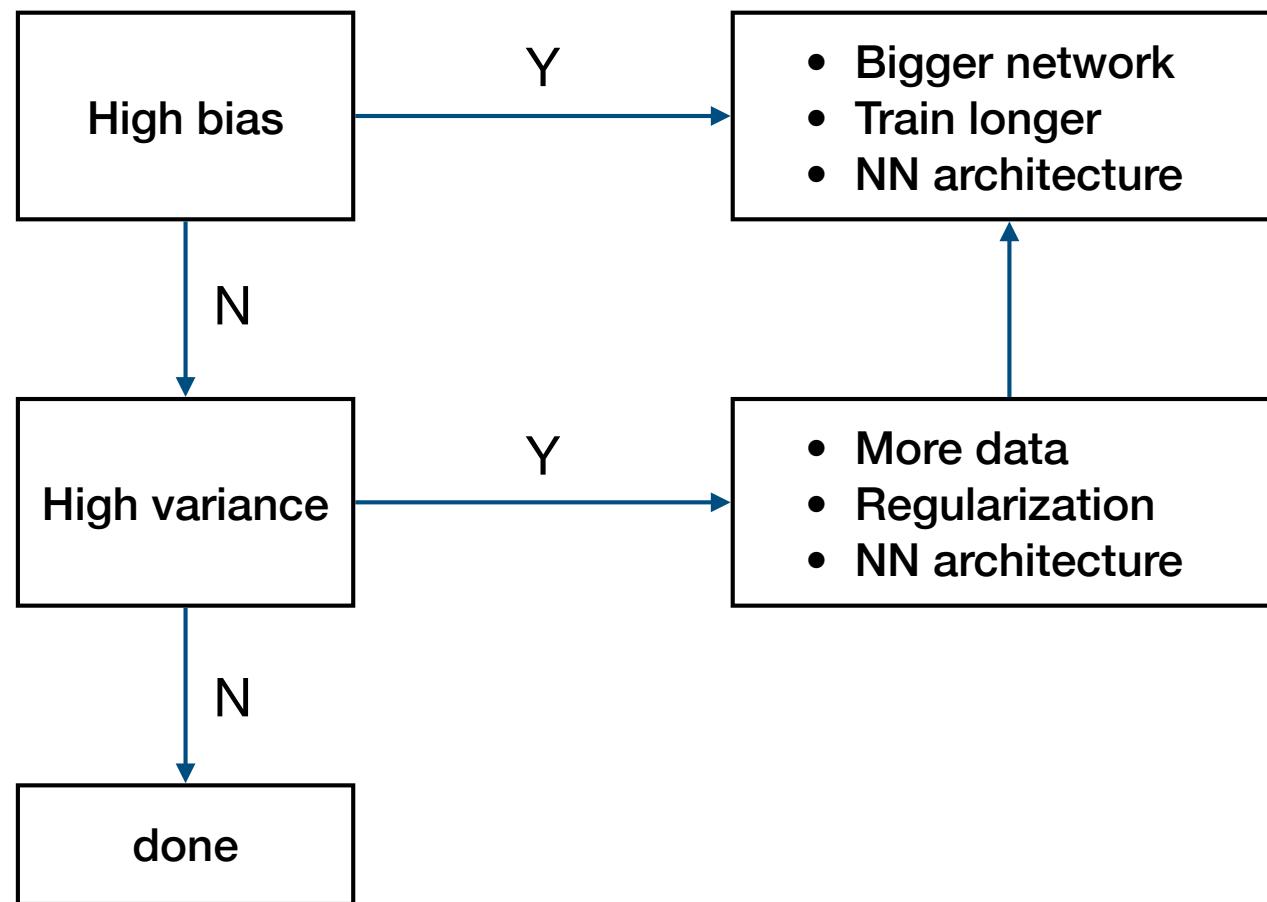
# Practical aspects of DNNs

- Bias and variance

	High variance	High bias	High bias High variance	Low Bias Low variance
Train set	1%	15%	15%	0.5%
Dev set	11%	16%	30%	1%

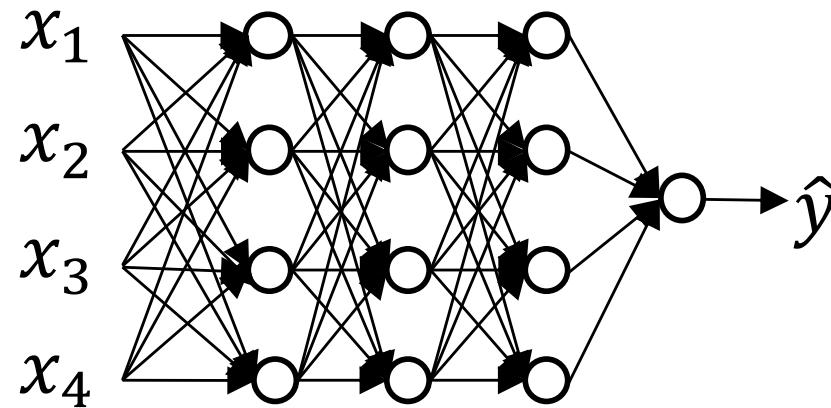


# Practical aspects of DNNs

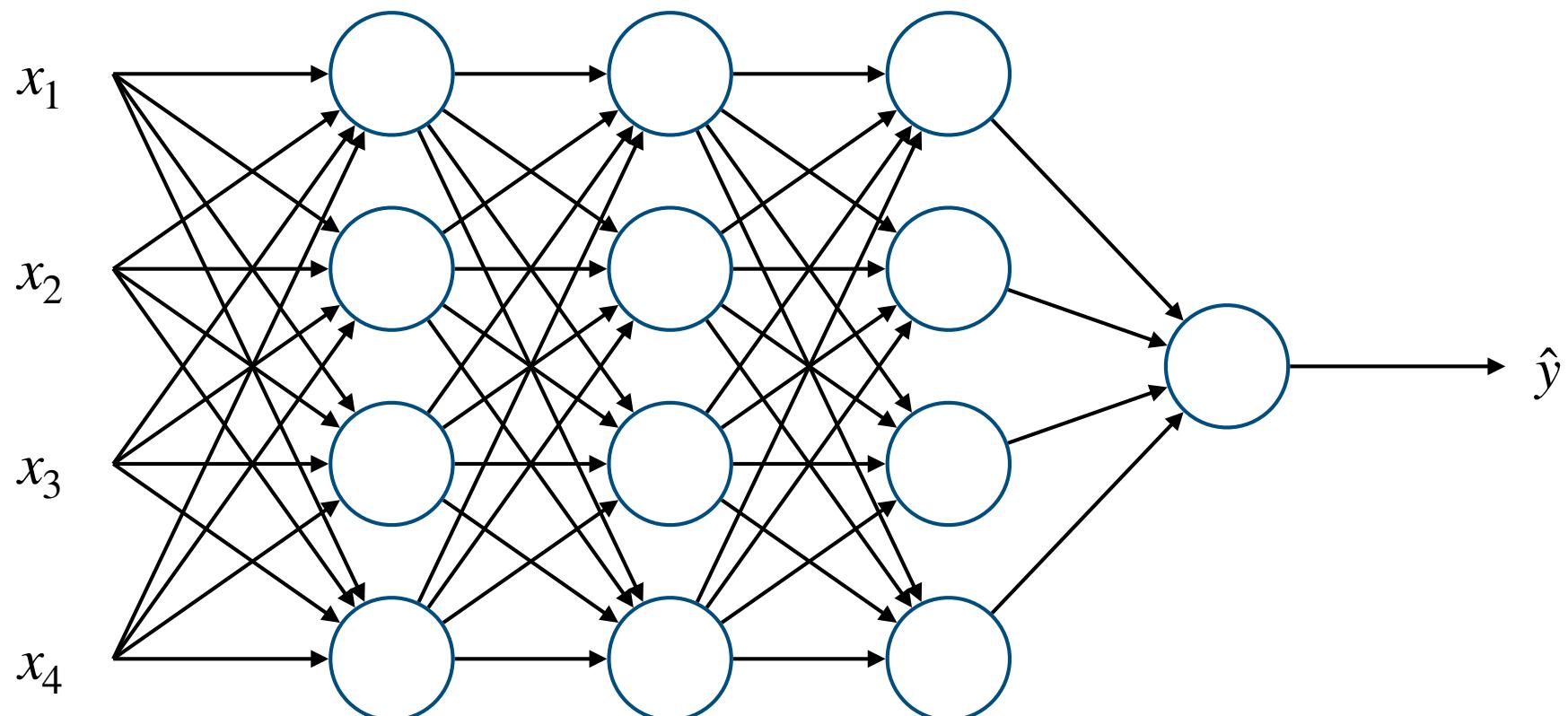


# Practical aspects of DNNs

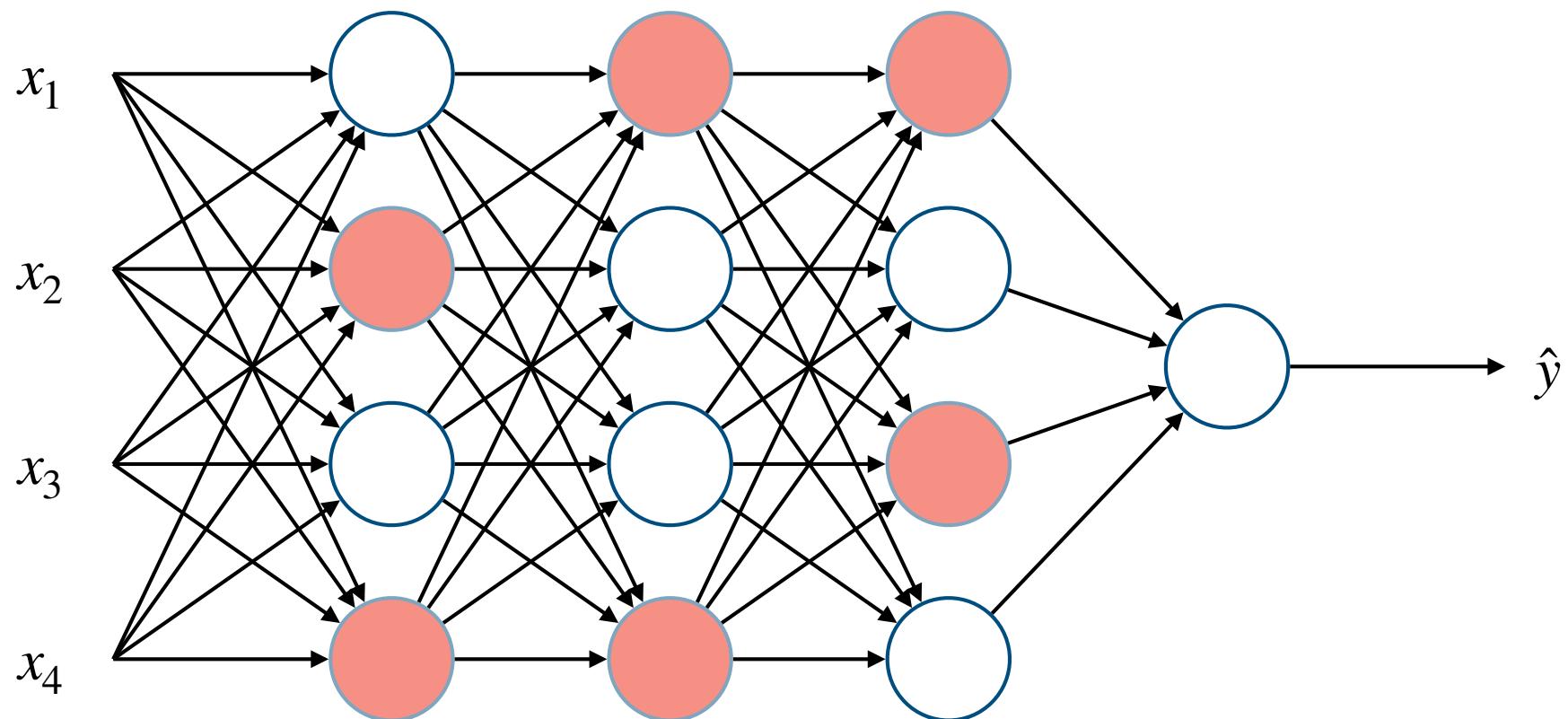
- Regularization can reduce overfitting
- Dropout regularization
- Intuition: cannot rely on any one feature, so have to spread out weights



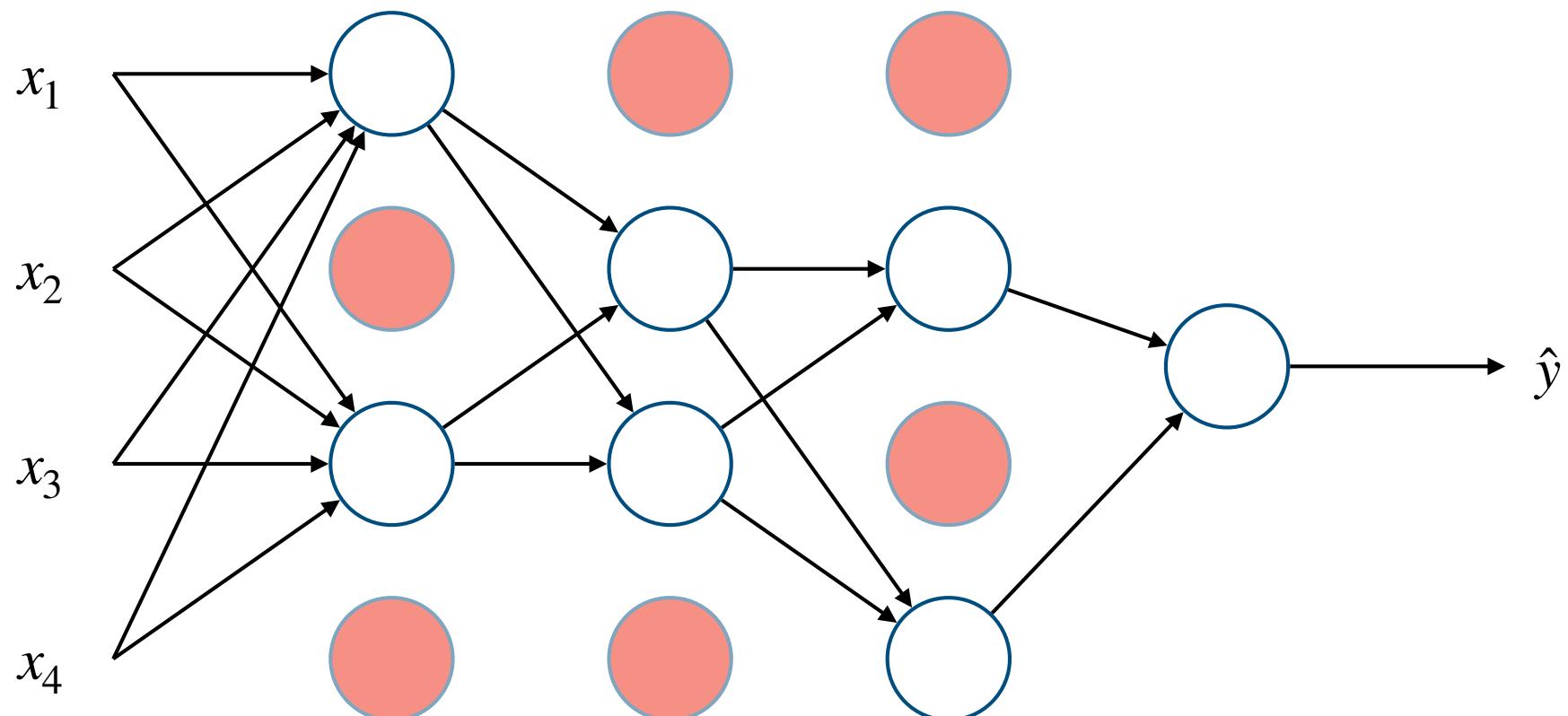
# Practical aspects of DNNs



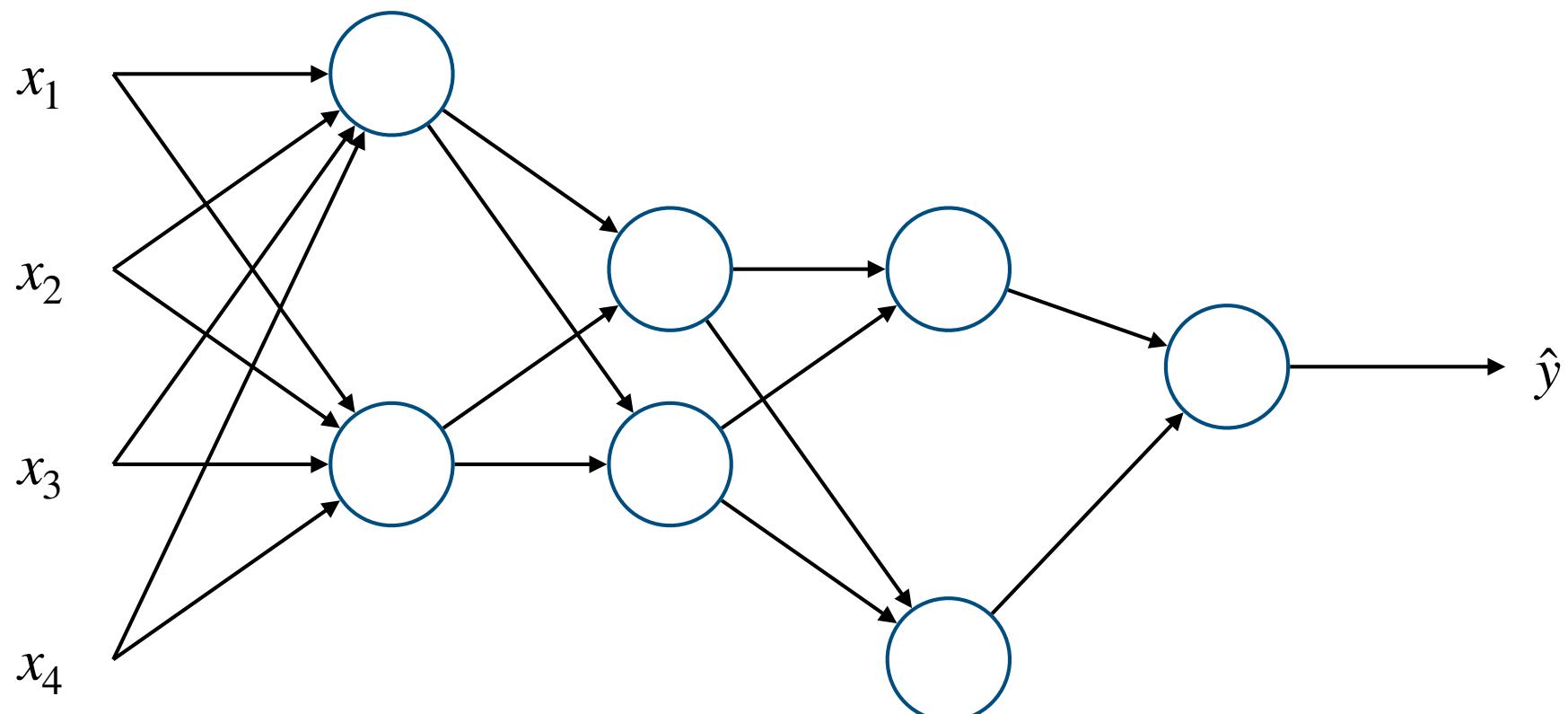
# Practical aspects of DNNs



# Practical aspects of DNNs



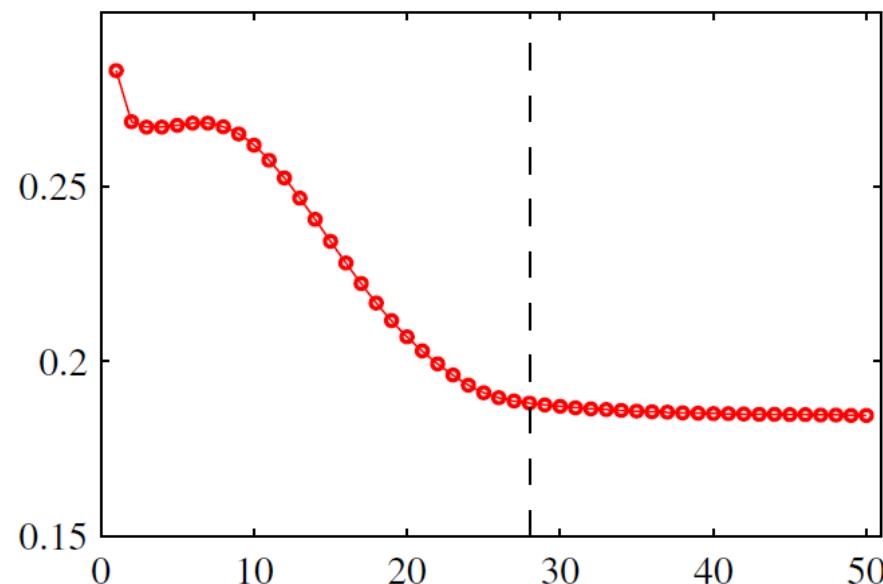
# Practical aspects of DNNs



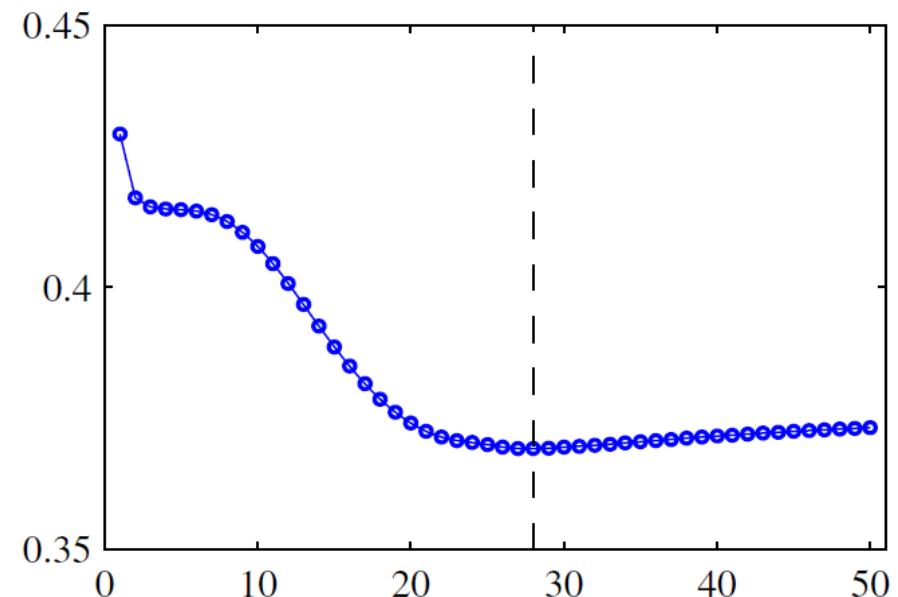
# Practical aspects of DNNs

- Alternative to regularization: early stopping
- Stop at smallest error with validation data

training set error



validation set error



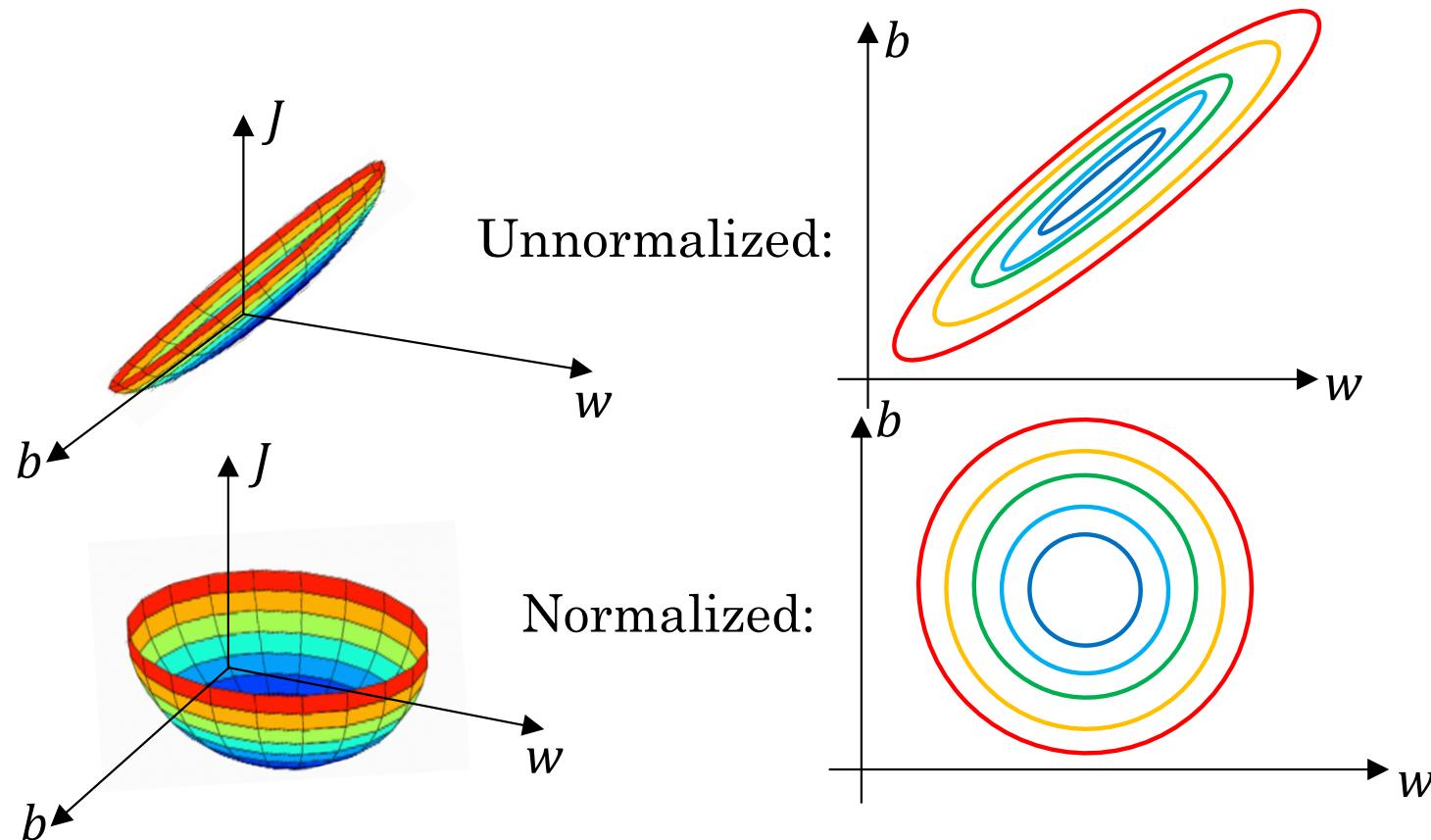
# Practical aspects of DNNs

- Data augmentation



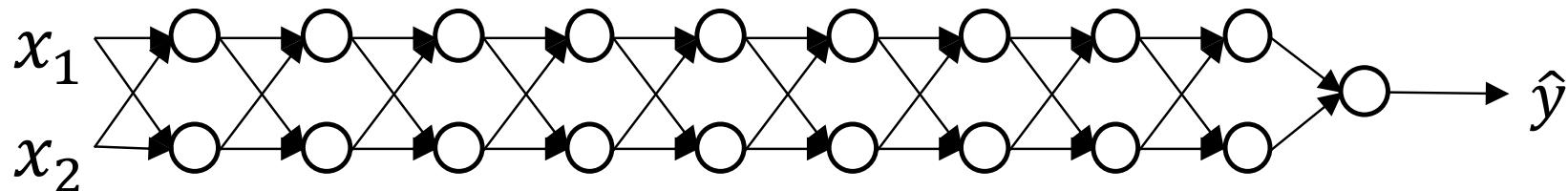
# Practical aspects of DNNs

- Normalize inputs



# Practical aspects of DNNs

- Vanishing / exploding gradients



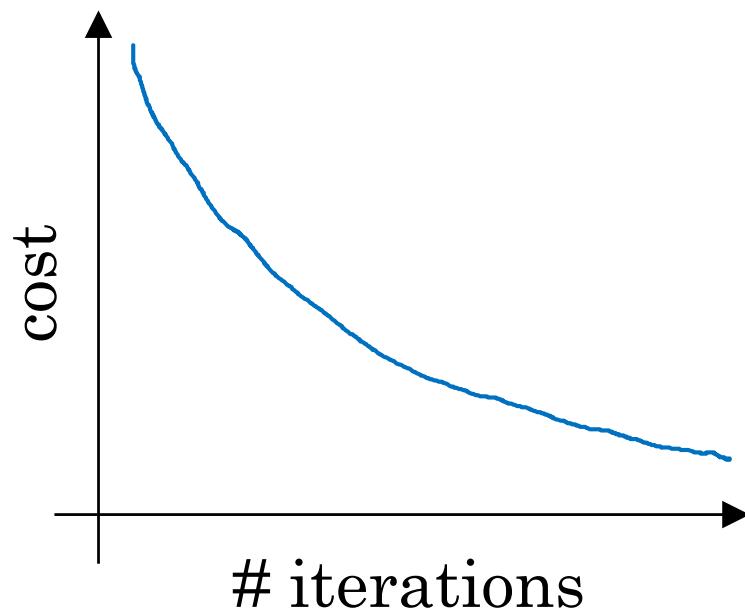
- Use linear activation function
- $\hat{y} = W^{[1]}W^{[2]} \dots W^{[l-1]}W^{[l]}x$

# Optimization algorithms

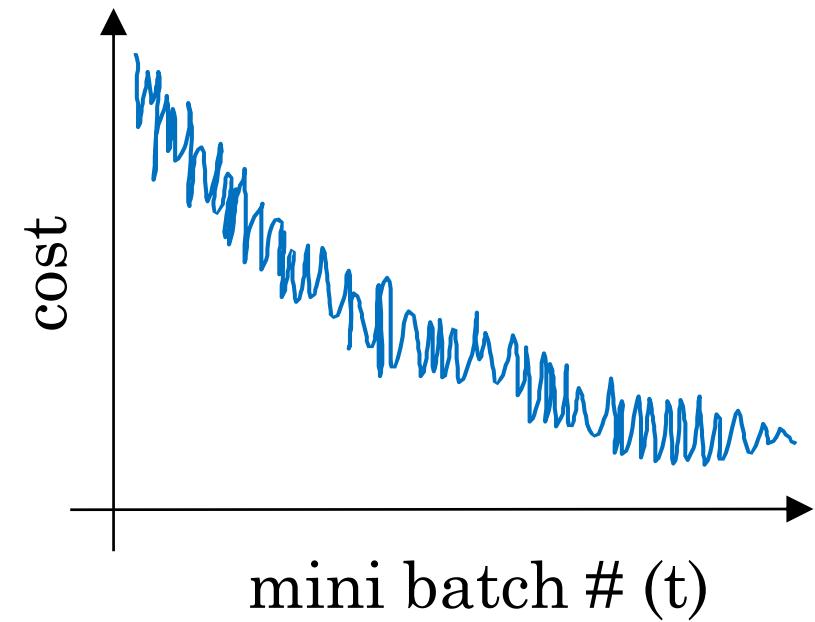
- Mini-batch gradient descent
  - > Vectorization allows efficient computation on many examples
  - > Forward propagation uses all data
  - > Idea: split up training set in batches and apply gradient descent on evaluation of one batch

# Optimization algorithms

Batch gradient descent

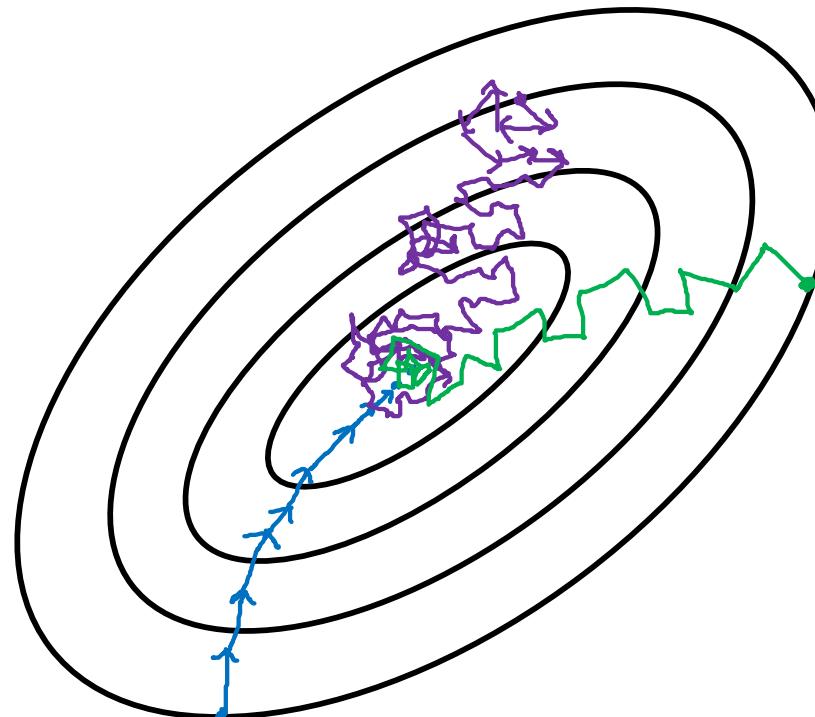


Mini-batch gradient descent



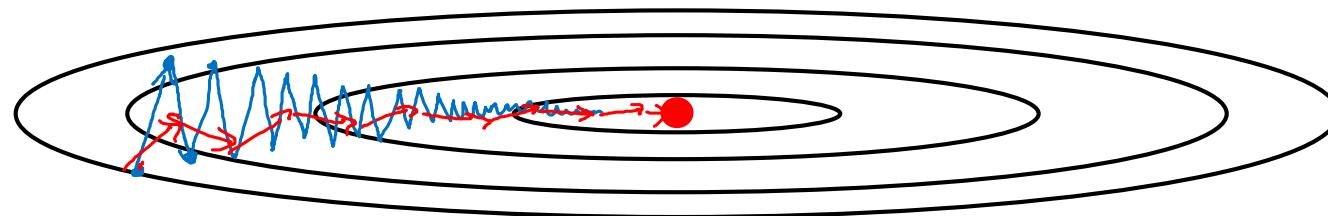
# Optimization algorithms

- Purple = stochastic gradient descent
- Green = mini-batch gradient descent
- Blue = batch gradient descent



# Optimization algorithms

- Gradient descent with momentum



- In practice, you want to dampen directions in which you should not move

# Optimization algorithms

- Gradient descent with momentum
- On iteration  $t$ , compute  $dW$ ,  $db$  on the current mini-batch

$$v_{dW} = \beta v_{dW} + (1 - \beta) dW$$

$$v_{db} = \beta v_{db} + (1 - \beta) db$$

$$W = W - \alpha v_{dW}$$

$$b = b - \alpha v_{db}$$

- Hyperparameters:  $\alpha, \beta$   $(\beta = 0.9)$

# Optimization algorithms

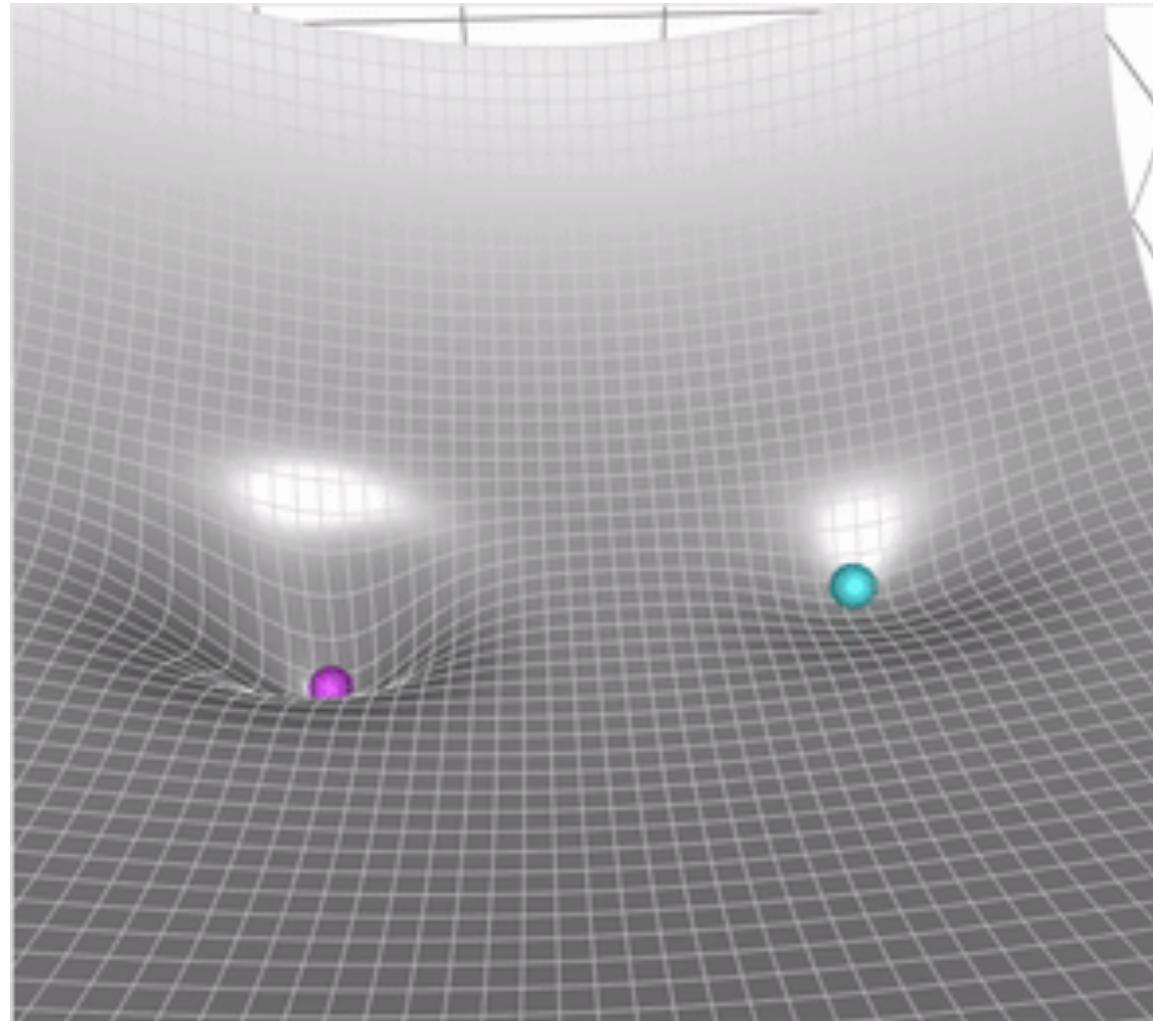
- Taking average is an example of smoothing

$$\bar{x}_{n+1} = \frac{x_1 + \cdots + x_{n+1}}{n+1} = \frac{x_1 + \cdots + x_n}{n+1} + \frac{x_{n+1}}{n+1}$$

$$= \frac{x_1 + \cdots + x_n}{n} \cdot \frac{n}{n+1} + \frac{x_{n+1}}{n+1} = \frac{n}{n+1} \bar{x}_n + \frac{1}{n+1} x_{n+1}$$

$$= \alpha_n \bar{x}_n + (1 - \alpha_n) x_{n+1}, \text{ with } \alpha_n = \frac{n}{n+1}$$

# Optimization algorithms



- gradient descent (cyan)
- momentum (magenta)

# Optimization algorithms

- RMSprop
- On iteration  $t$ , compute  $\mathbf{d}W$ ,  $\mathbf{db}$  on the current mini-batch

$$s_{\mathbf{d}W} = \beta s_{\mathbf{d}W} + (1 - \beta) \mathbf{d}W^2$$

$$s_{\mathbf{d}b} = \beta s_{\mathbf{d}b} + (1 - \beta) \mathbf{d}b^2$$

$$W = W - \alpha \frac{\mathbf{d}W}{\sqrt{s_{\mathbf{d}W}}}$$

$$b = b - \alpha \frac{\mathbf{d}b}{\sqrt{s_{\mathbf{d}b}}}$$

# Optimization algorithms

- Properties of RMSprop

$$s_{\mathbf{d}W} = \beta s_{\mathbf{d}W} + (1 - \beta) \mathbf{d}W^2 \rightarrow s_t = \beta s_{t-1} + (1 - \beta) g_t^2$$

- Note that  $s_{t-1} = \beta s_{t-2} + (1 - \beta) g_{t-1}^2$ . By substitution, we get

$$s_t = \beta [\beta s_{t-2} + (1 - \beta) g_{t-1}^2] + (1 - \beta) g_t^2 = \beta^2 s_{t-2} + (1 - \beta)[g_t^2 + \beta g_{t-1}^2]$$

- By continuing this process, we get

$$s_t = (1 - \beta) \sum_{i=1}^t \beta^{t-i} g_i^2$$

# Optimization algorithms

- If we assume that the  $g_i$ 's are i.i.d., then

$$\mathbb{E}[s_t] = \mathbb{E} \left[ (1 - \beta) \sum_{i=1}^t \beta^{t-i} g_i^2 \right] = (1 - \beta) \sum_{i=1}^t \beta^{t-i} \mathbb{E}[g_i^2]$$

$$\mathbb{E}[s_t] = \mathbb{E}[g_t^2](1 - \beta) \sum_{i=1}^t \beta^{t-i} = \mathbb{E}[g_t^2](1 - \beta) \sum_{i=0}^{t-1} \beta^i$$

$$\mathbb{E}[s_t] = \mathbb{E}[g_t^2](1 - \beta) \frac{1 - \beta^t}{1 - \beta} = \mathbb{E}[g_t^2](1 - \beta^t)$$

-

# Optimization algorithms

- Adam optimization
- On iteration  $t$ , compute  $\mathbf{d}W$ ,  $\mathbf{db}$  on the current mini-batch

$$v_{\mathbf{d}W} = \beta_1 v_{\mathbf{d}W} + (1 - \beta_1) \mathbf{d}W$$

$$s_{\mathbf{d}W} = \beta_2 s_{\mathbf{d}W} + (1 - \beta_2) \mathbf{d}W^2$$

$$v_{\mathbf{d}b} = \beta_1 v_{\mathbf{d}b} + (1 - \beta_1) \mathbf{d}b$$

$$s_{\mathbf{d}b} = \beta_2 s_{\mathbf{d}b} + (1 - \beta_2) \mathbf{d}b^2$$

$$V_{\mathbf{d}W}^{\text{cor}} = v_{\mathbf{d}W} / (1 - \beta_1^t)$$

$$S_{\mathbf{d}W}^{\text{cor}} = s_{\mathbf{d}W} / (1 - \beta_2^t)$$

$$V_{\mathbf{d}b}^{\text{cor}} = v_{\mathbf{d}b} / (1 - \beta_1^t)$$

$$S_{\mathbf{d}b}^{\text{cor}} = s_{\mathbf{d}b} / (1 - \beta_2^t)$$

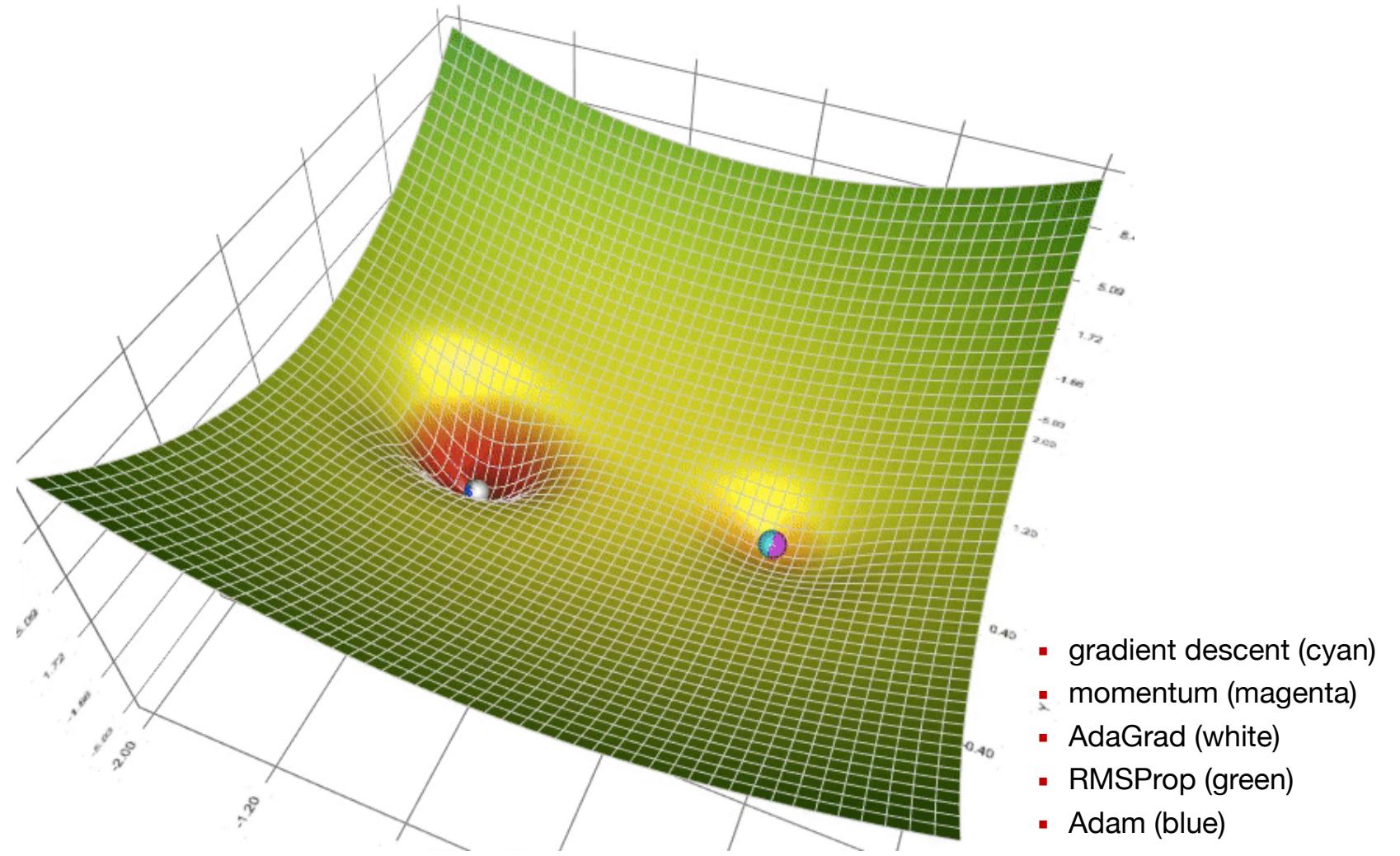
# Optimization algorithms

- Adam optimization

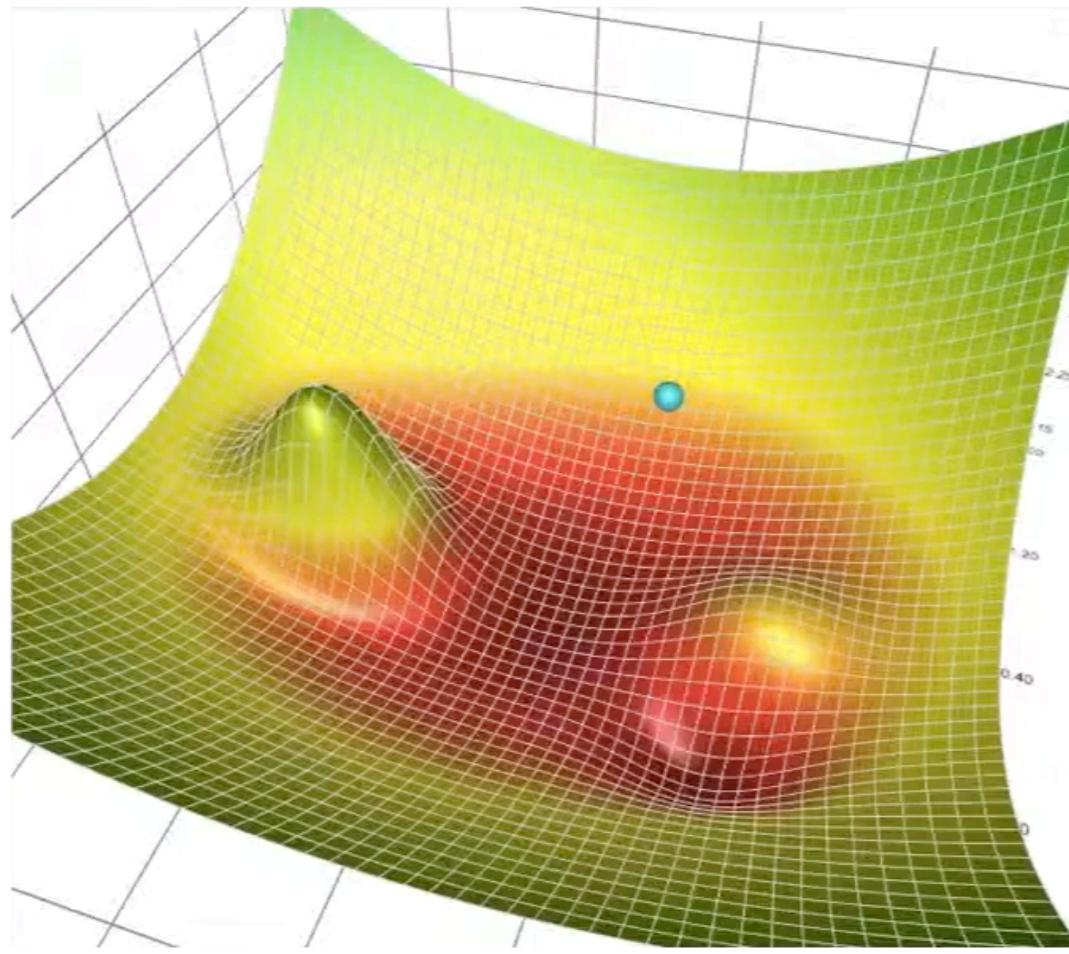
$$W = W - \alpha \frac{V_{dw}^{\text{cor}}}{\sqrt{S_{dw}^{\text{cor}}} + \epsilon}$$

$$b = b - \alpha \frac{V_{db}^{\text{cor}}}{\sqrt{S_{db}^{\text{cor}}} + \epsilon}$$

# Optimization algorithms

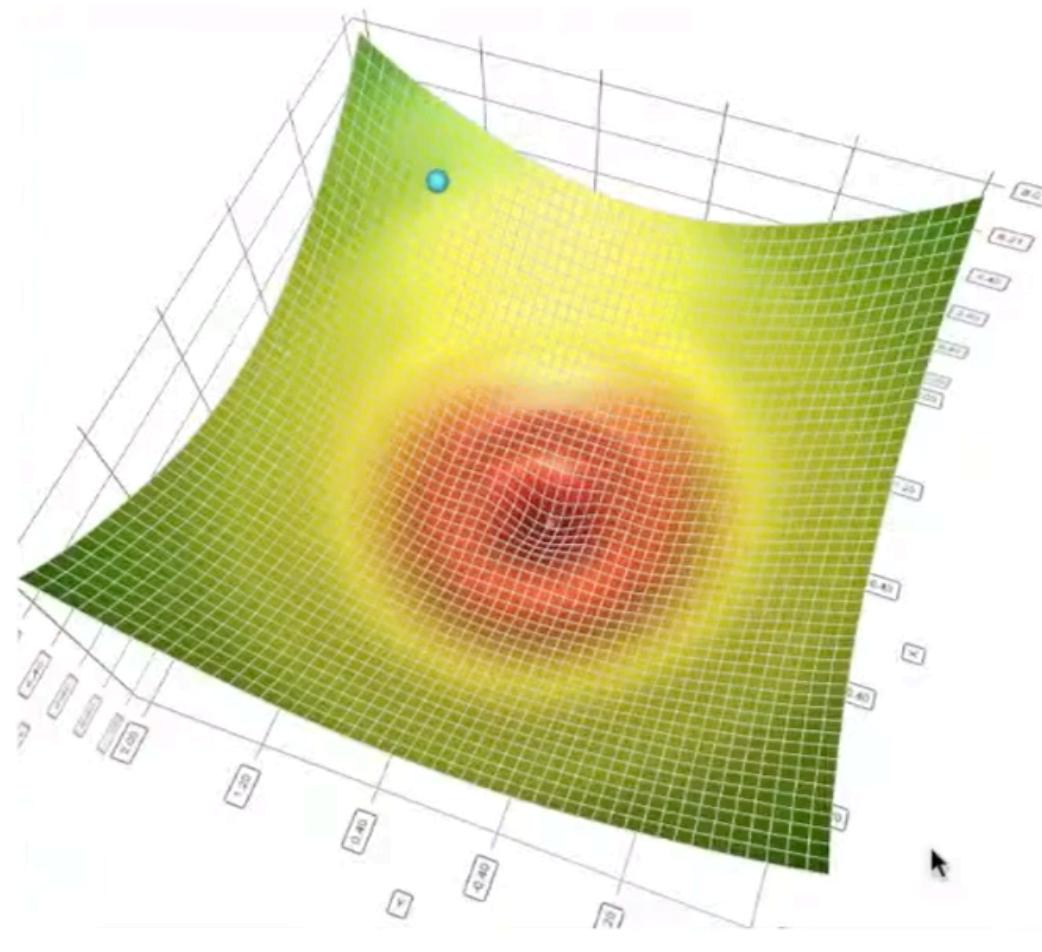


# Optimization algorithms



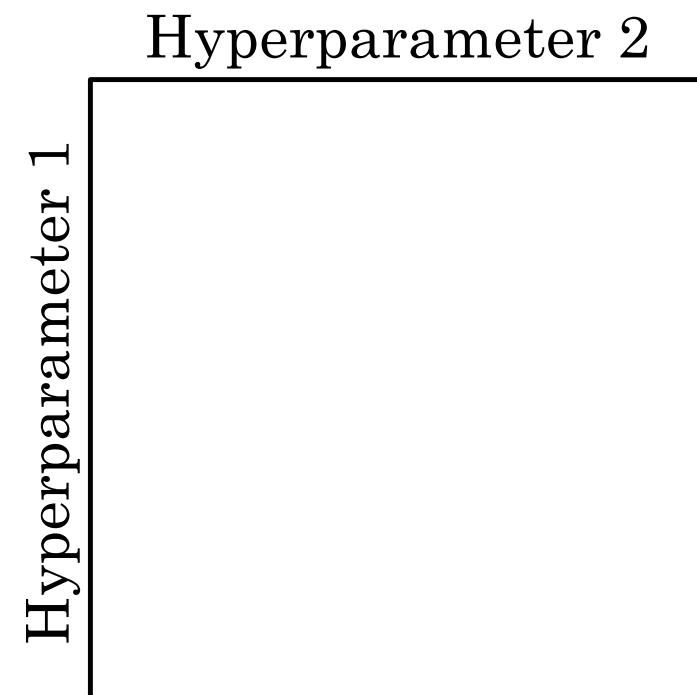
- gradient descent (cyan)
- momentum (magenta)
- AdaGrad (white)
- RMSProp (green)
- Adam (blue)

# Optimization algorithms

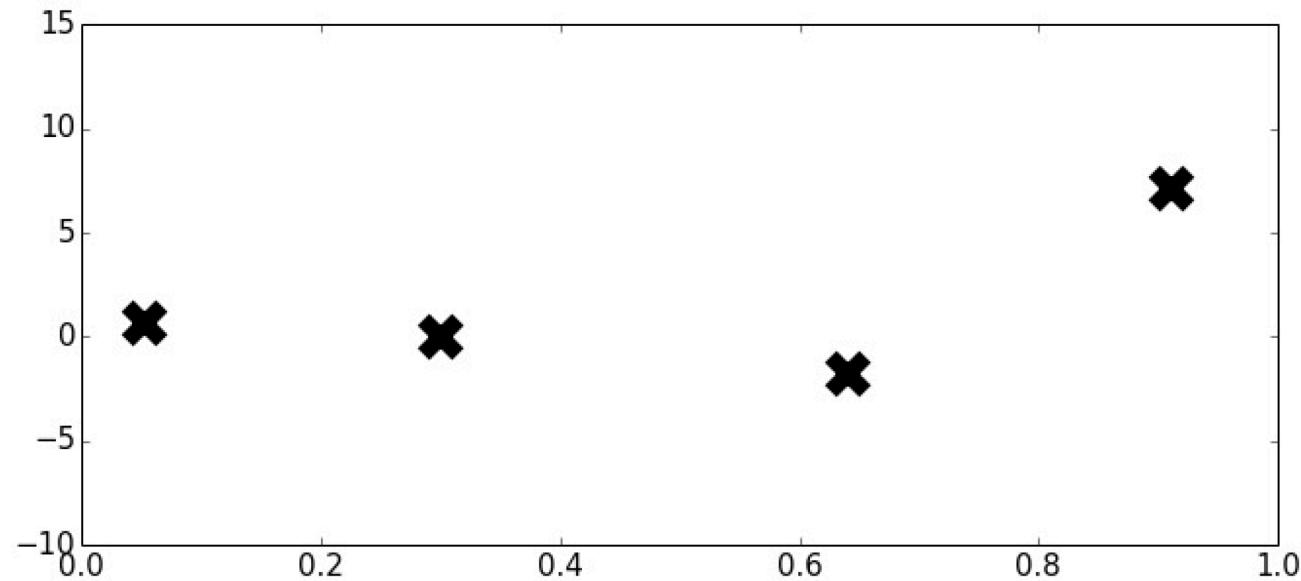


# Optimization algorithms

- Hyperparameter tuning:
  - Do not use a grid
  - Go from coarse to fine
  - Pick appropriate scale
  - Run in parallel



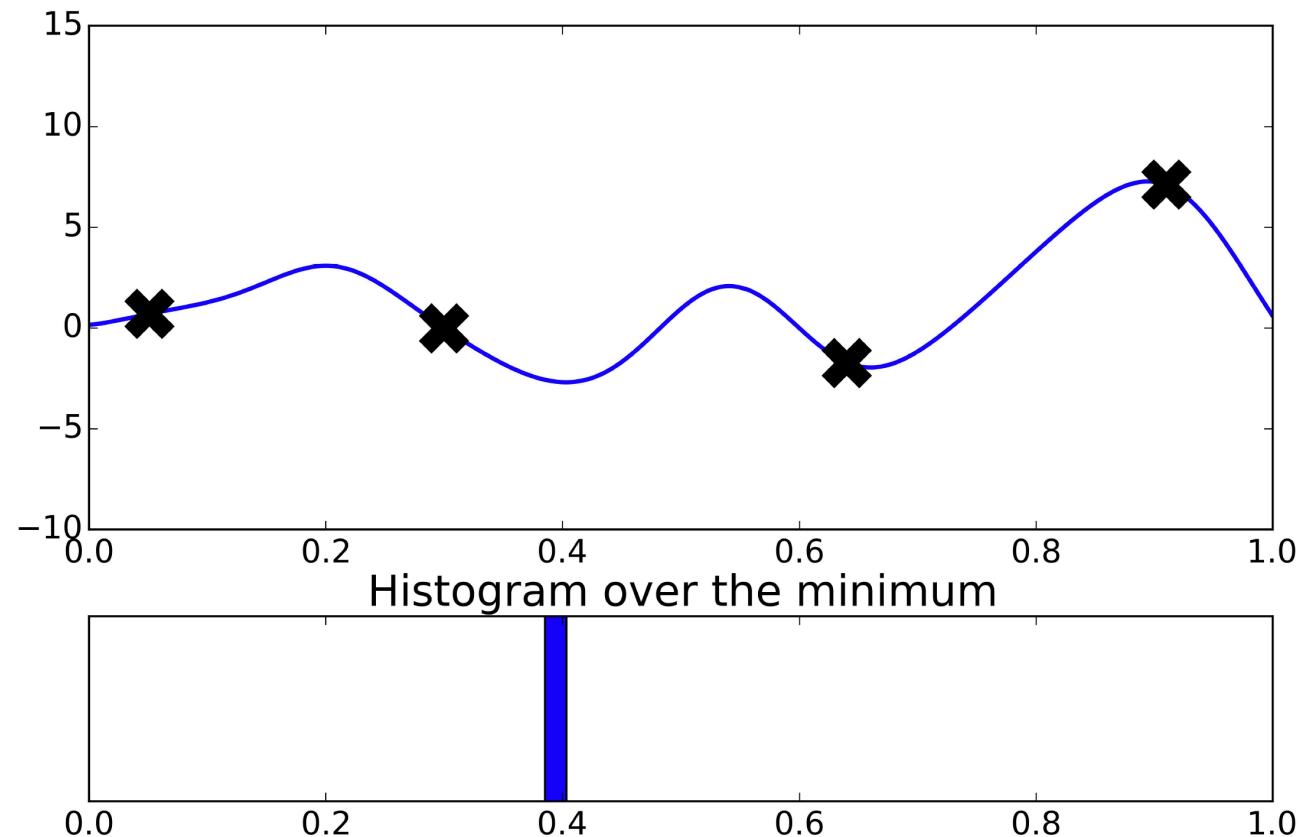
# Optimization algorithms



- Where is the minimum of  $f$ ?
- Where should the next evaluation be taken?

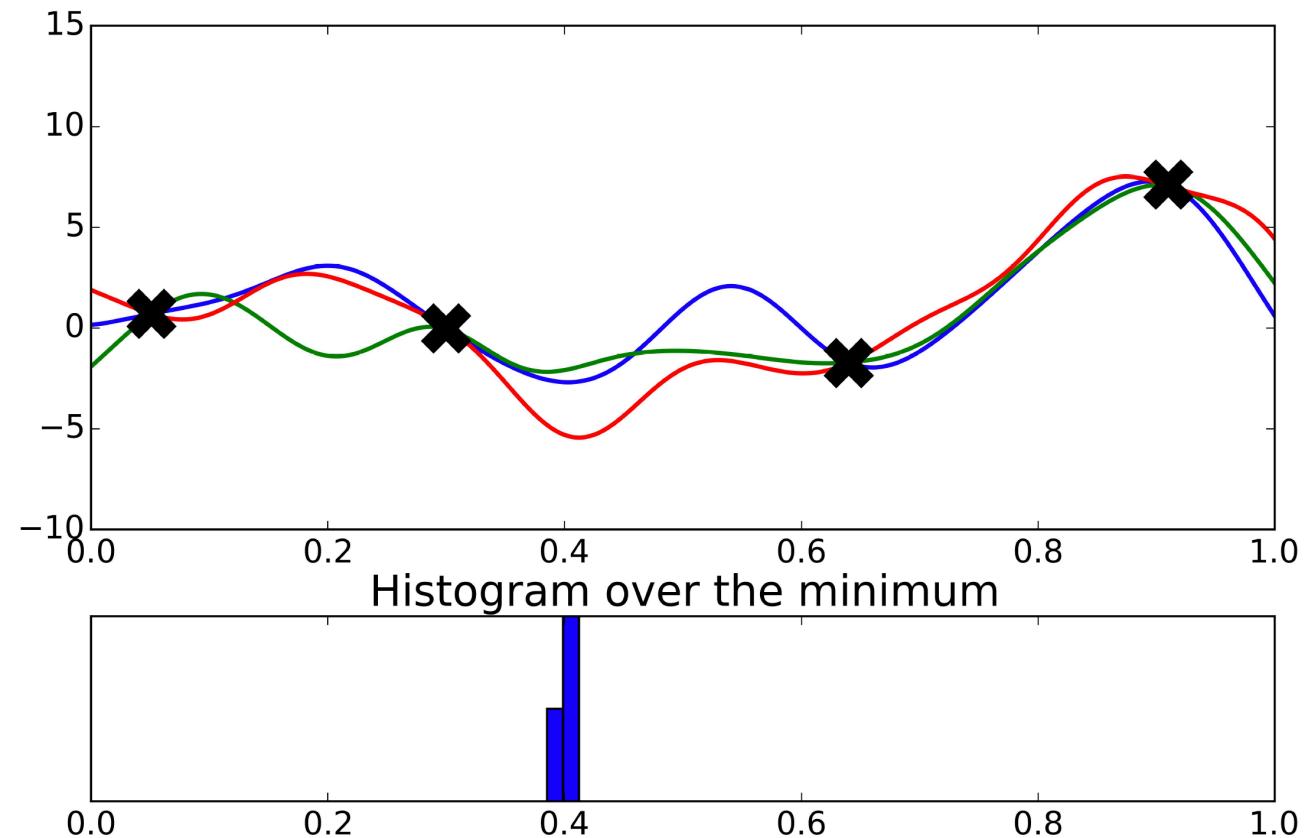
# Optimization algorithms

- Hyperparameter tuning: Bayesian optimization – 1 curve



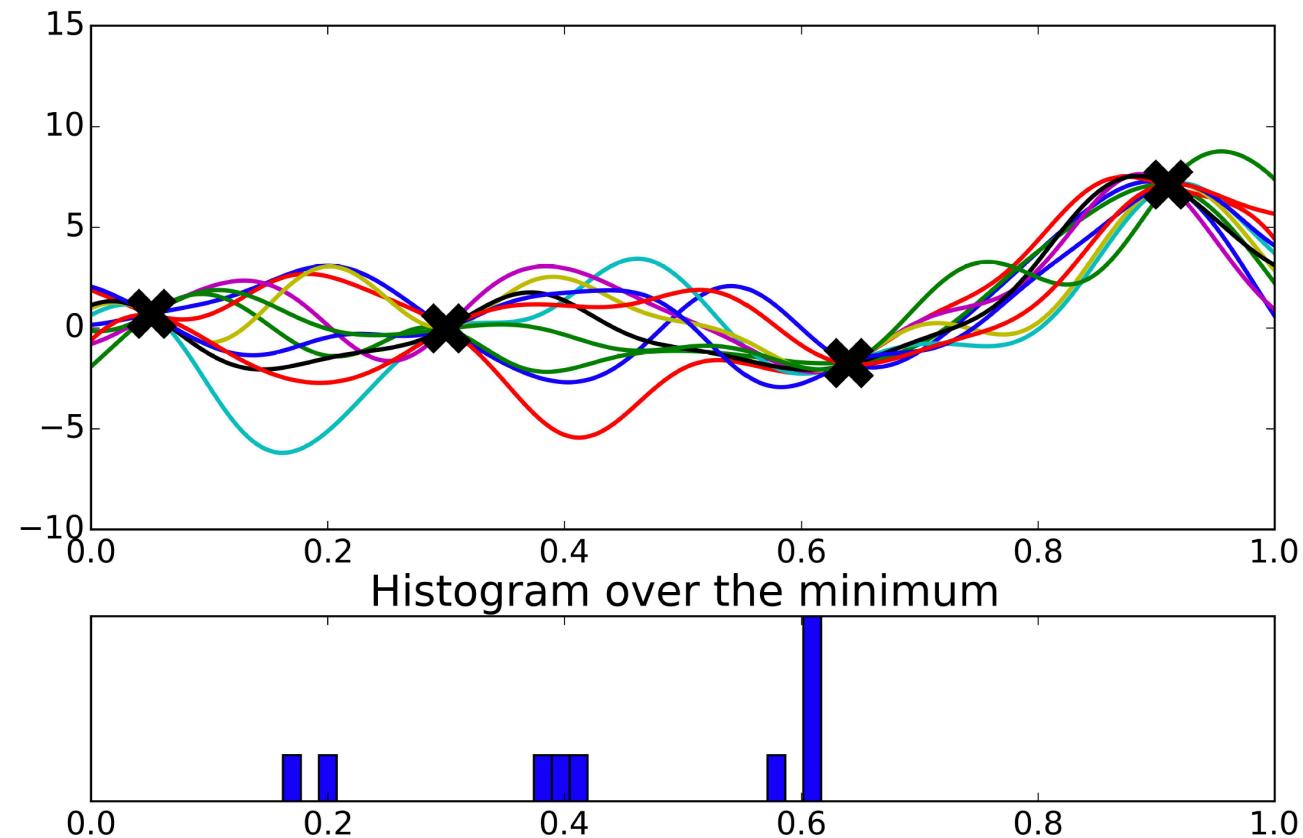
# Optimization algorithms

- Hyperparameter tuning: Bayesian optimization – 3 curves



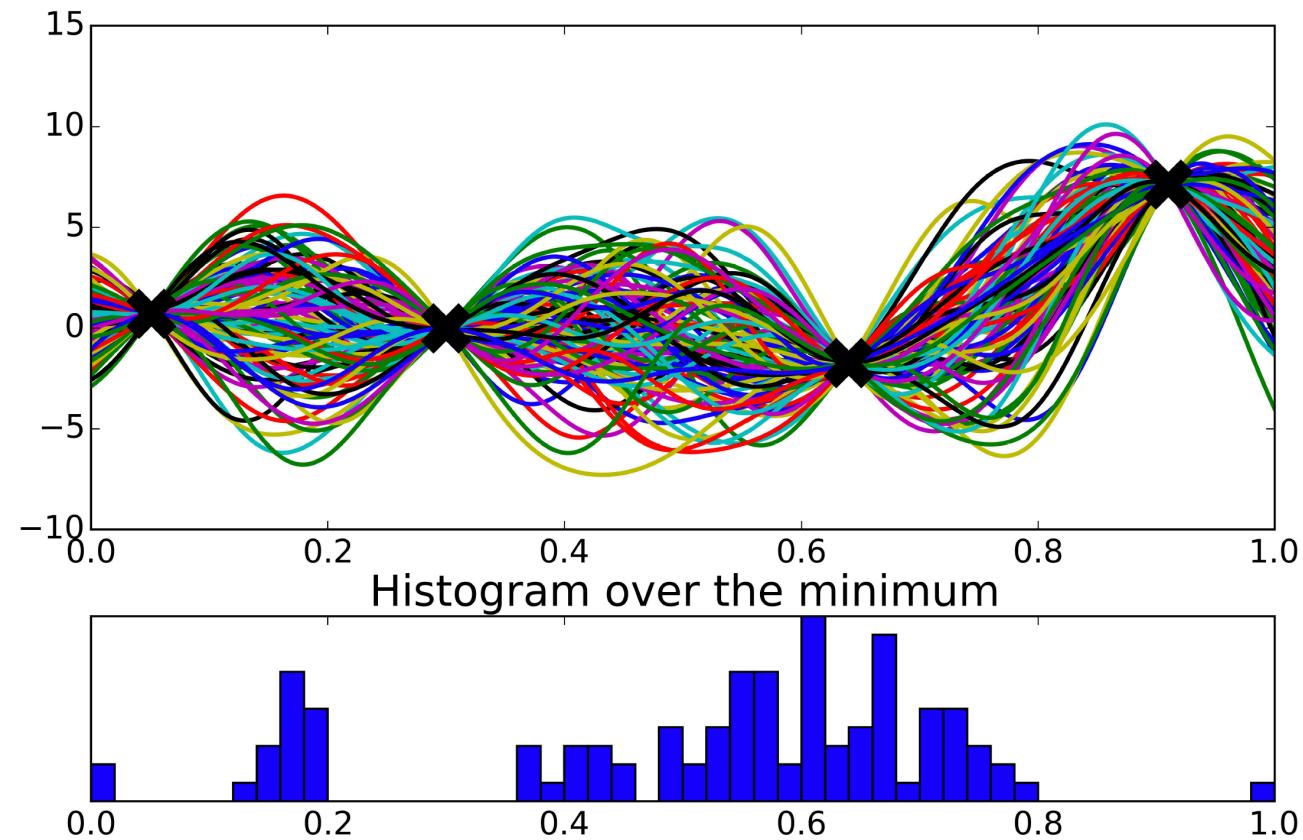
# Optimization algorithms

- Hyperparameter tuning: Bayesian optimization – 10 curves



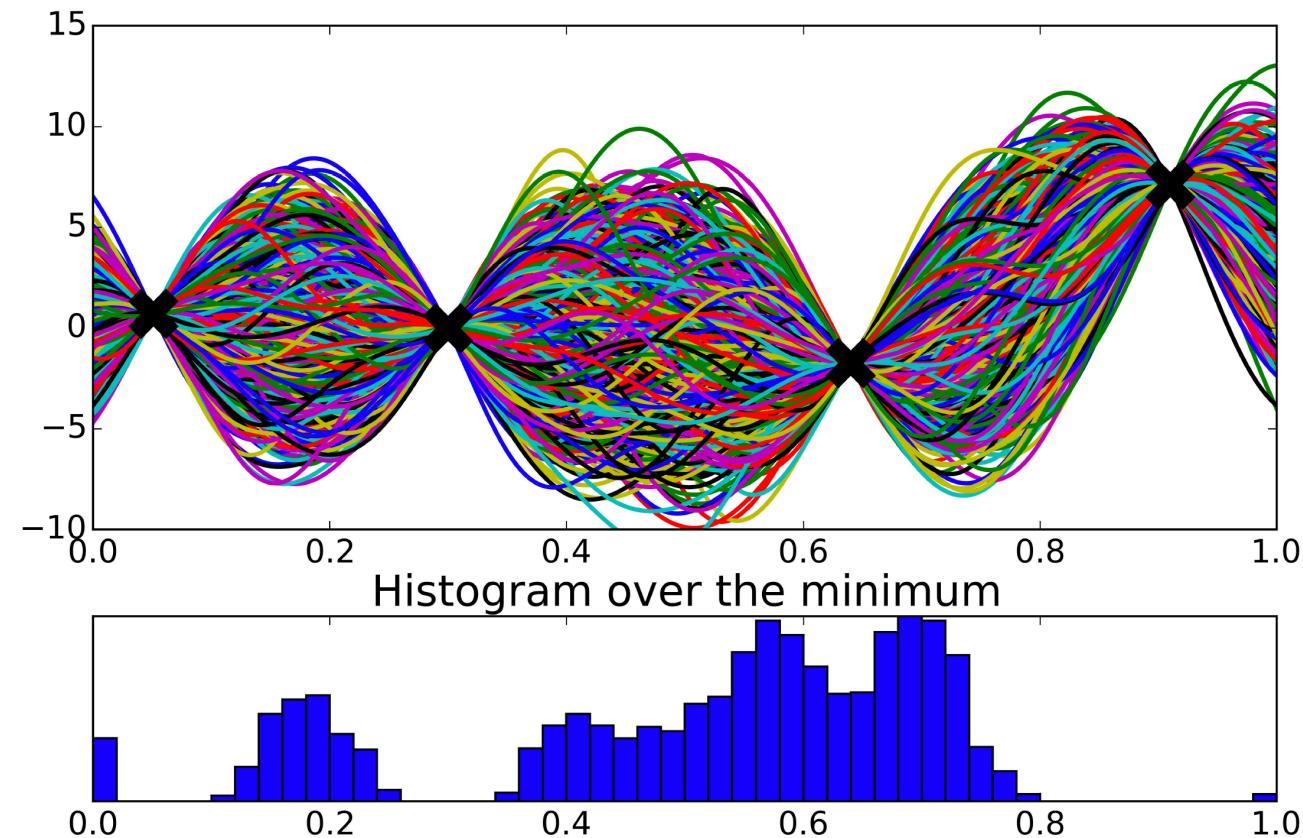
# Optimization algorithms

- Hyperparameter tuning: Bayesian optimization – 100 curves



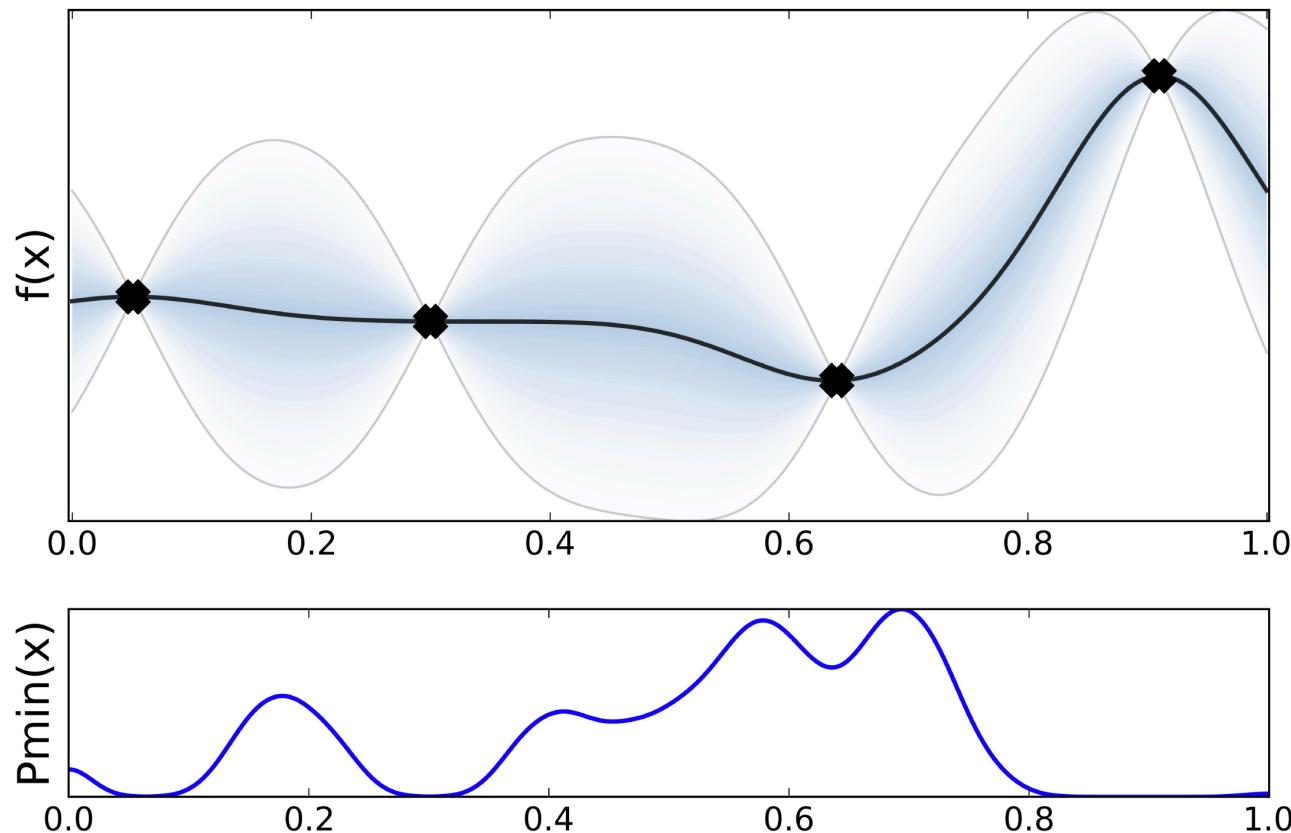
# Optimization algorithms

- Hyperparameter tuning: Bayesian optimization – many curves



# Optimization algorithms

- Hyperparameter tuning: Bayesian optimization –  $\infty$  curves

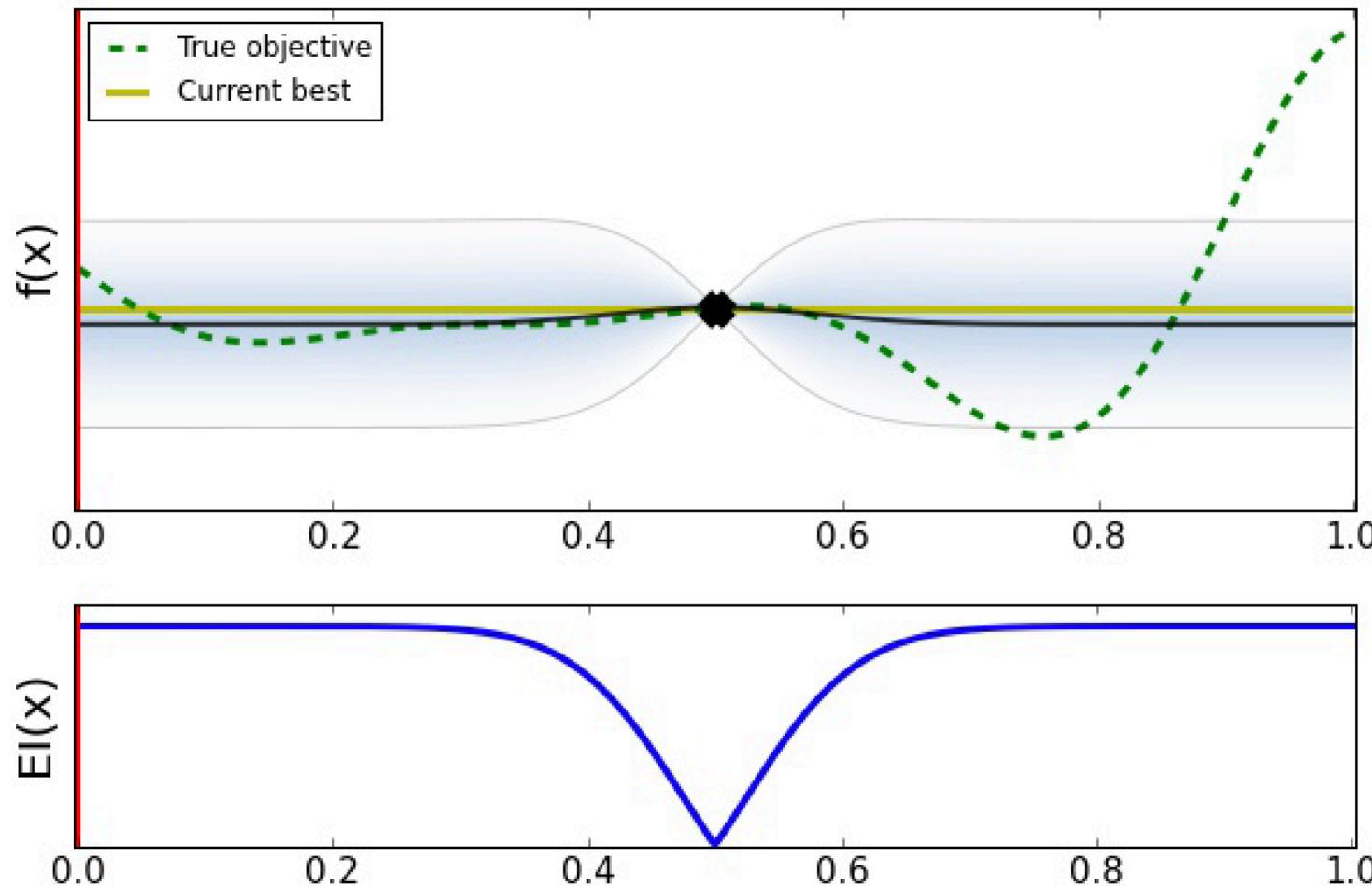


# Optimization algorithms

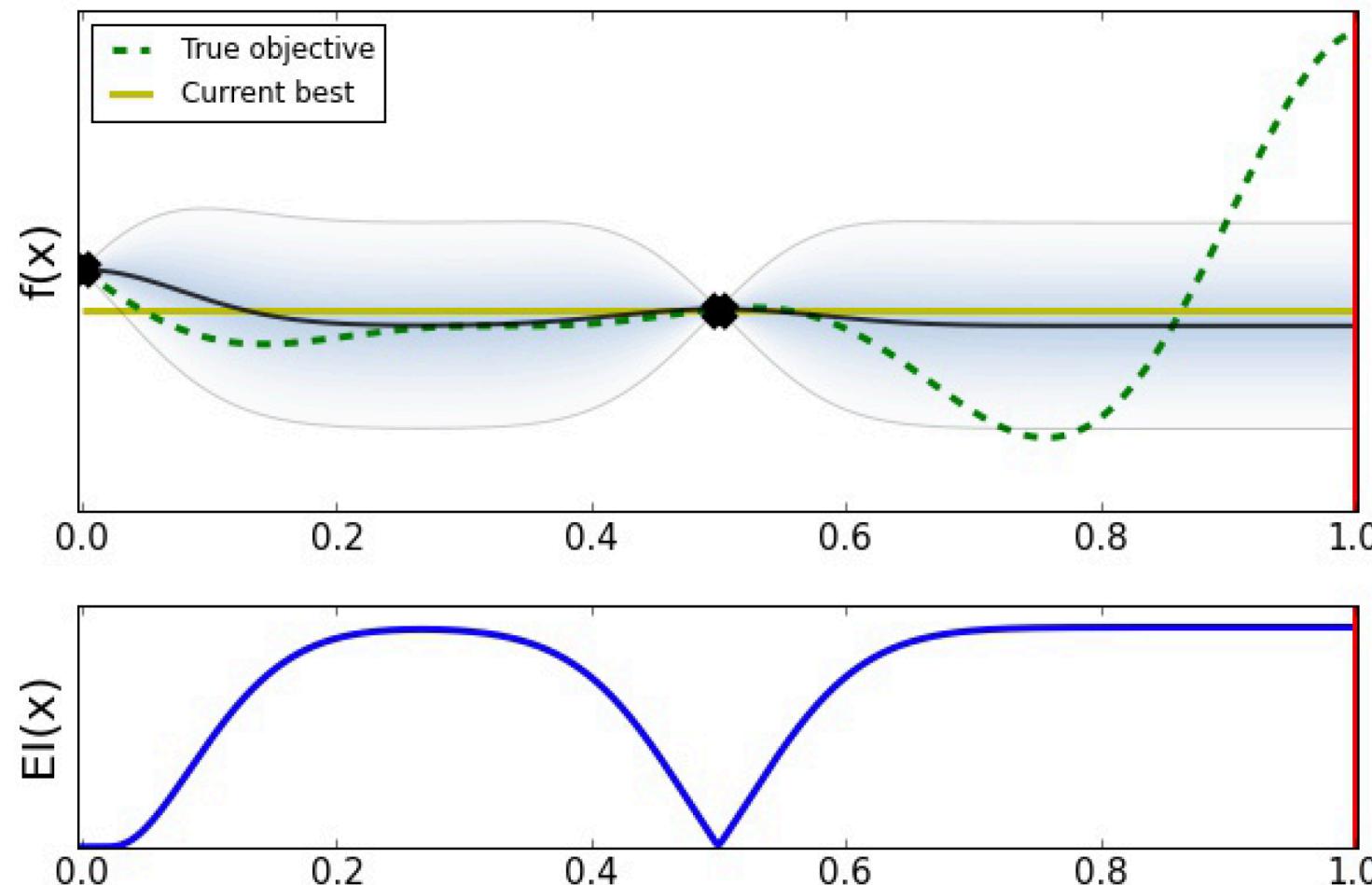
- Use a surrogate model of  $f$  to carry out the optimization
  - >  $f(x) \sim \text{GP}(\mu(x), k(x, x'))$
  - >  $\mu(x)$  is the mean function
  - >  $k(x, x'; \theta)$  is the covariance function
- Expected improvement  $EI(x)$  is given by

$$EI(\mathbf{x}; \theta, \mathcal{D}) = \int_y \max(0, y_{\text{best}} - y) p(y | \mathbf{x}; \theta, \mathcal{D}) dy$$

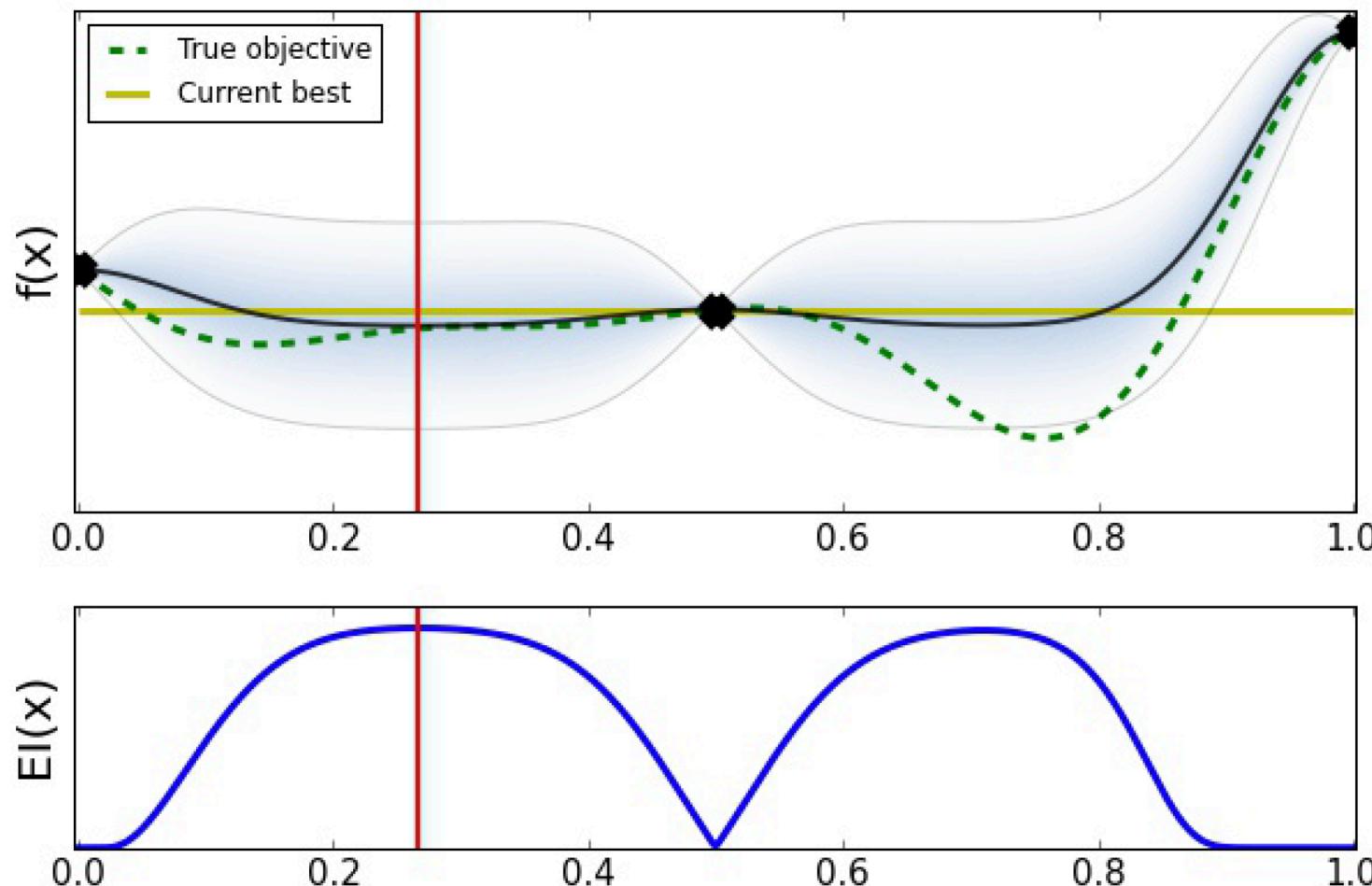
# Optimization algorithms



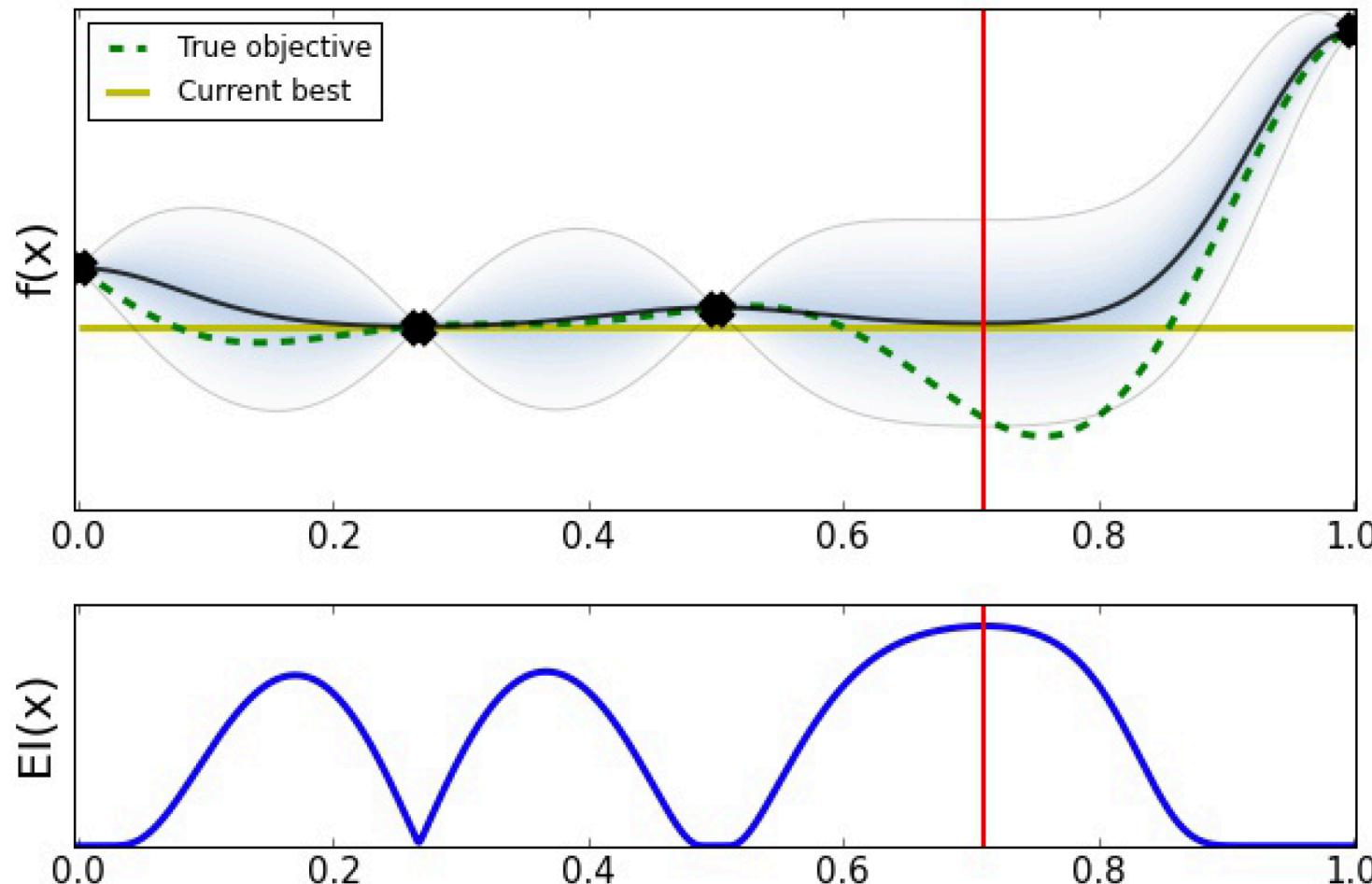
# Optimization algorithms



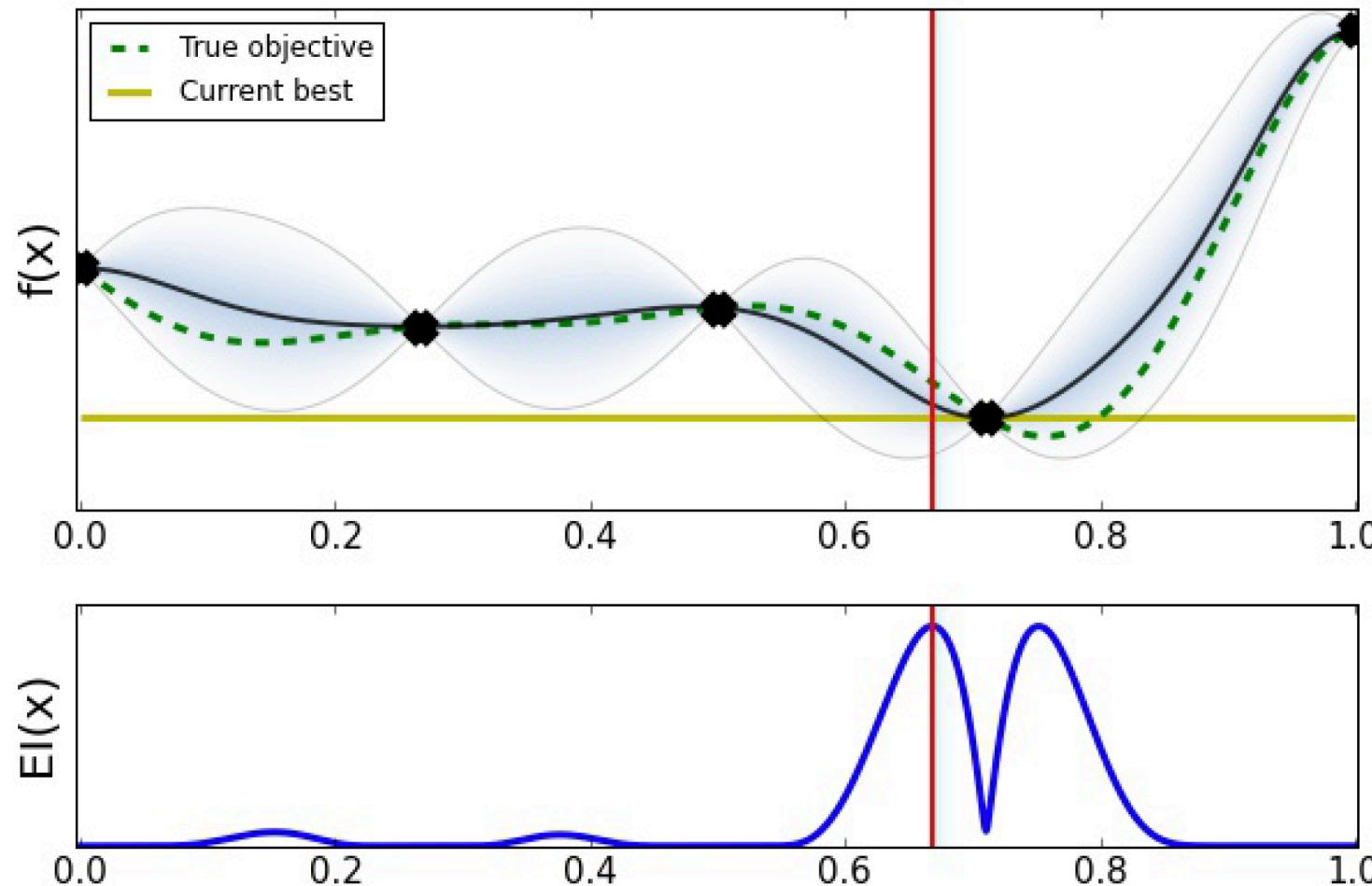
# Optimization algorithms



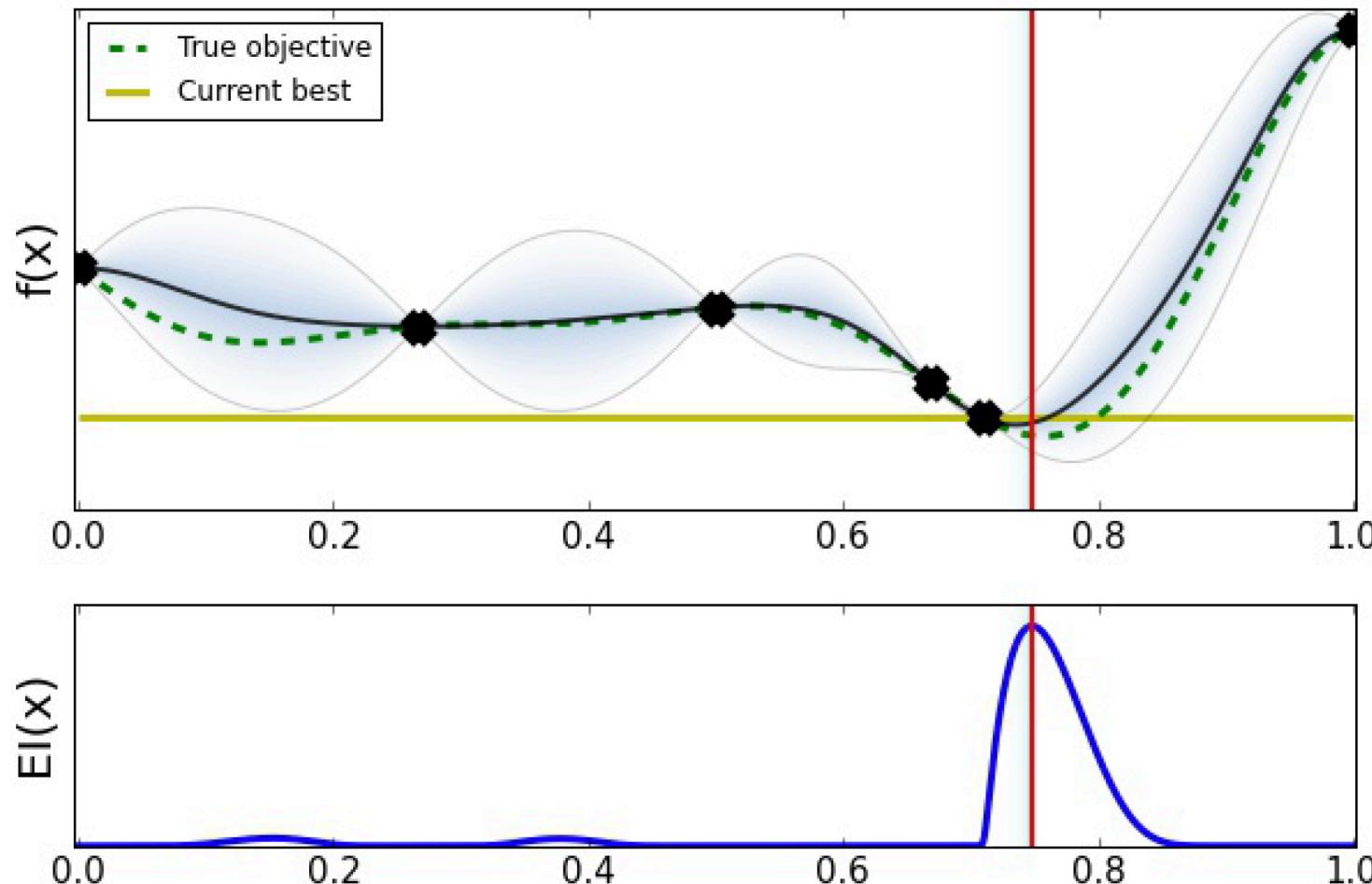
# Optimization algorithms



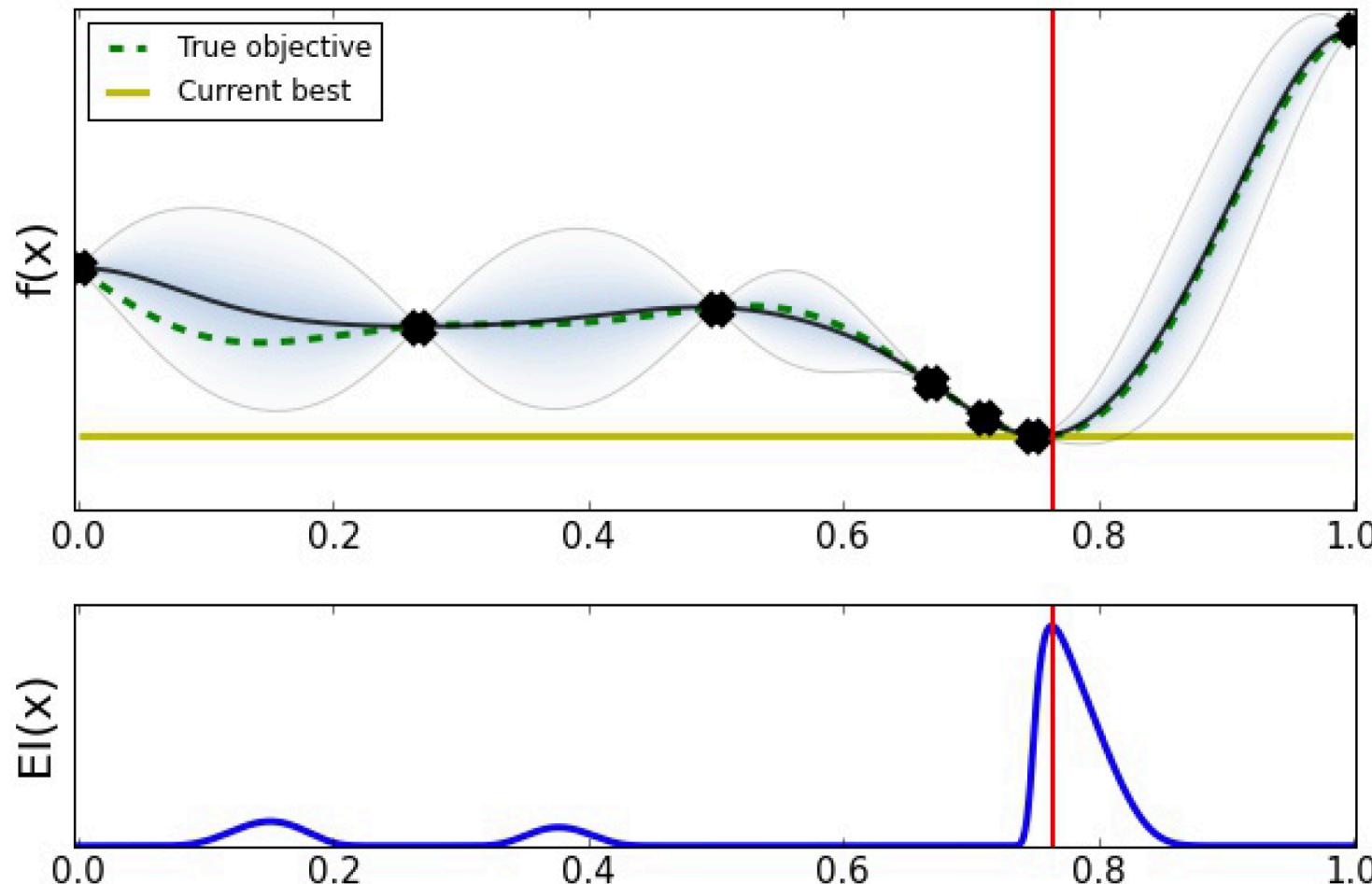
# Optimization algorithms



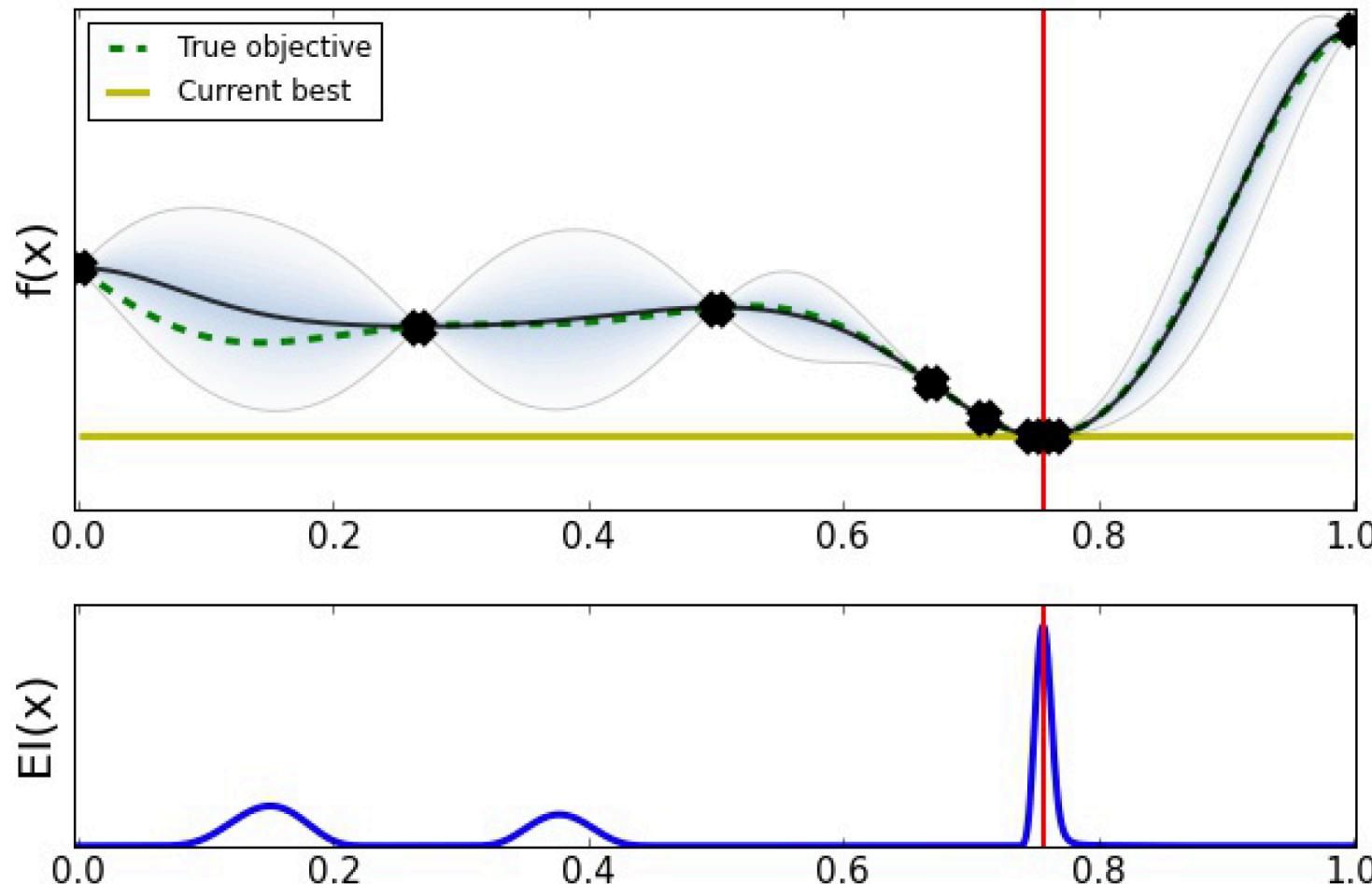
# Optimization algorithms



# Optimization algorithms



# Optimization algorithms



# Optimization algorithms

- Caffe/Caffe2
- CNTK
- DL4J
- Keras
- Lasagne
- mxnet
- PaddlePaddle
- TensorFlow
- Theano
- Torch