

Lecture Notes for **Machine Learning in Python**

Professor Eric Larson
Adaptive Neural Network Optimization

Class Logistics and Agenda

- Agenda:
 - More optimization techniques
 - Review
 - Adaptive Learning Strategies
 - Town Hall for MLP

Class Overview, by topic

Table Data
Visualization

Numpy, Pandas, Seaborn
Overviews with some in-depth discussion

Dimension
Reduction and
Image Processing

Scikit-learn, Scikit Image,
Intuition only, Some mathematics

Linear and
Logistic
Regression

Numpy, Recreate API for Scikit-learn
Detailed mathematics for simple optimization
intuition for advanced optimization

Neural Networks
and Back Prop.

Numpy
Detailed mathematics for NN operations

Wide and Deep
Networks

Convolutional
Networks

Recurrent
Networks

Keras, Tensorflow
Intuition, Detailed implement.

Ethics in
Language Models

ConceptNet
Case studies

$$\mathbf{W}_{k+1} = \mathbf{W}_k - \rho_k$$

$$\mathbf{V}^{(2)} = \mathbf{A}^{(3)} - \mathbf{Y}$$

new final layer update

- Cross entropy

- Momentum

$$\rho_k = \alpha \nabla J(\mathbf{W}_k) + \beta \nabla J(\mathbf{W}_{k-1})$$

- Nesterov's Accelerated Gradient

$$\rho_k = \underbrace{\beta \nabla J(\mathbf{W}_k + \alpha \nabla J(\mathbf{W}_{k-1}))}_{\text{step twice}} + \alpha \nabla J(\mathbf{W}_{k-1})$$

- Mini-batching

← all data →

| | batch 1 | batch 2 | batch 3 | batch 4 | batch 5 | batch 6 | batch 7 | batch 8 | batch 9 |
|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|
| Epoch 1 | | | | | | | | | |
| Epoch 2 | | | | | | | | | |
| Epoch 3 | | | | | | | | | |
| Epoch 4 | | | | | | | | | |
| ... | | | | | | | | | |

shuffle ordering each epoch and update W 's after each batch

- Learning rate adaptation (eta)

$$\eta_e = \eta_0^{(1+e \cdot \epsilon)}$$

Review: Activations Summary

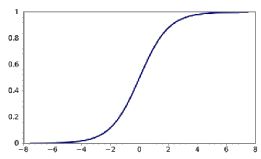
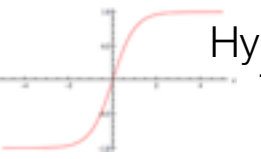
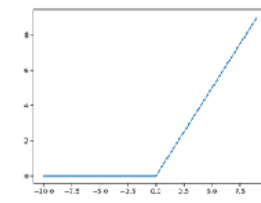
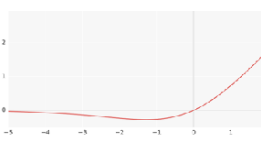
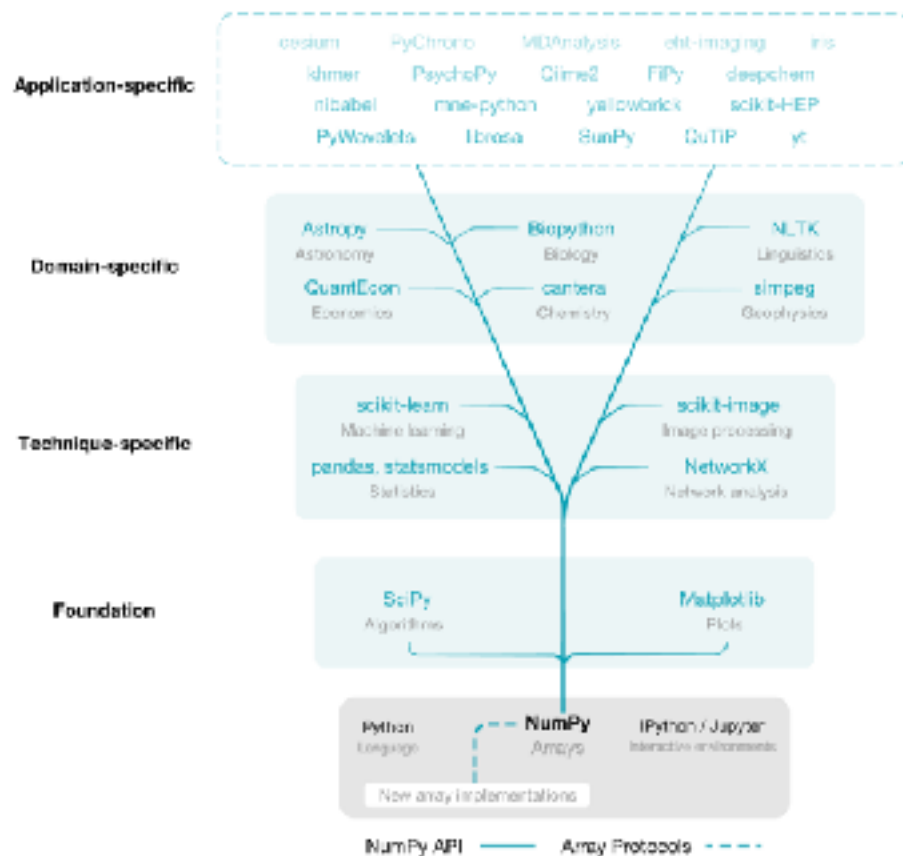
| | Definition | Derivative | Weight Init (Uniform Bounds) |
|--|---|--|--|
|  <p>Sigmoid</p> | $\phi(z) = \frac{1}{1 + e^{-z}}$ | $\nabla \phi(z) = a(1 - a)$ | $w_{ij}^{(L)} \approx \pm 4 \sqrt{\frac{6}{n^{(L)} + n^{(L+1)}}}$ |
|  <p>Hyperbolic Tangent</p> | $\phi(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$ | $\nabla \phi(z) = \frac{4}{(e^z + e^{-z})^2}$ | $w_{ij}^{(L)} \approx \pm \sqrt{\frac{6}{n^{(L)} + n^{(L+1)}}}$ |
|  <p>ReLU</p> | $\phi(z) = \begin{cases} z, & \text{if } z > 0 \\ 0, & \text{else} \end{cases}$ | $\nabla \phi(z) = \begin{cases} 1, & \text{if } z > 0 \\ 0, & \text{else} \end{cases}$ | $w_{ij}^{(L)} \approx \pm \sqrt{2} \sqrt{\frac{6}{n^{(L)} + n^{(L+1)}}}$ |
|  <p>SiLU</p> | $\phi(z) = \frac{z}{1 + e^{-z}}$ | $\nabla \phi(z) = a + \frac{(1 - a)}{1 + e^{-z}}$ | $w_{ij}^{(L)} \approx \pm \sqrt{2} \sqrt{\frac{6}{n^{(L)} + n^{(L+1)}}}$ |

Fig. 2: NumPy is the base of the scientific Python ecosystem.

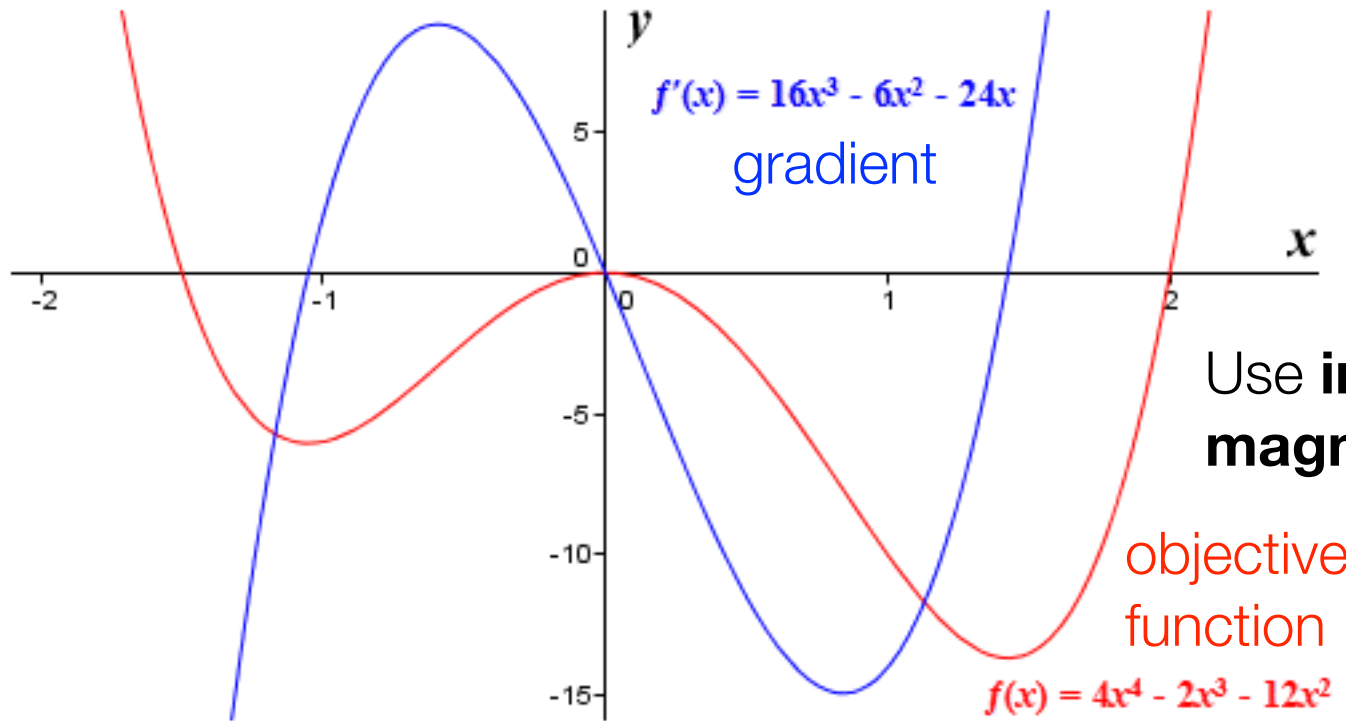
From: Array programming with NumPy

More Adaptive Optimization



Be adaptive based on Gradient Magnitude?

- Decelerate down regions that are steep
- Accelerate on plateaus



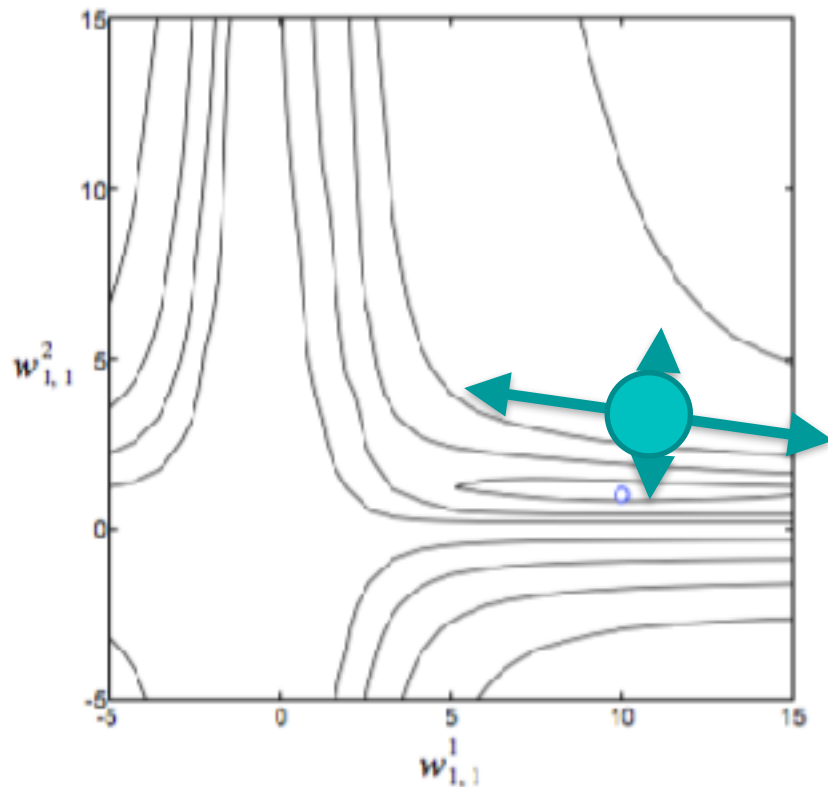
Use **inverse** of
magnitude of **gradient**!

Also **accumulate inverse** to be robust to
abrupt changes in **steepness**... momentum!!

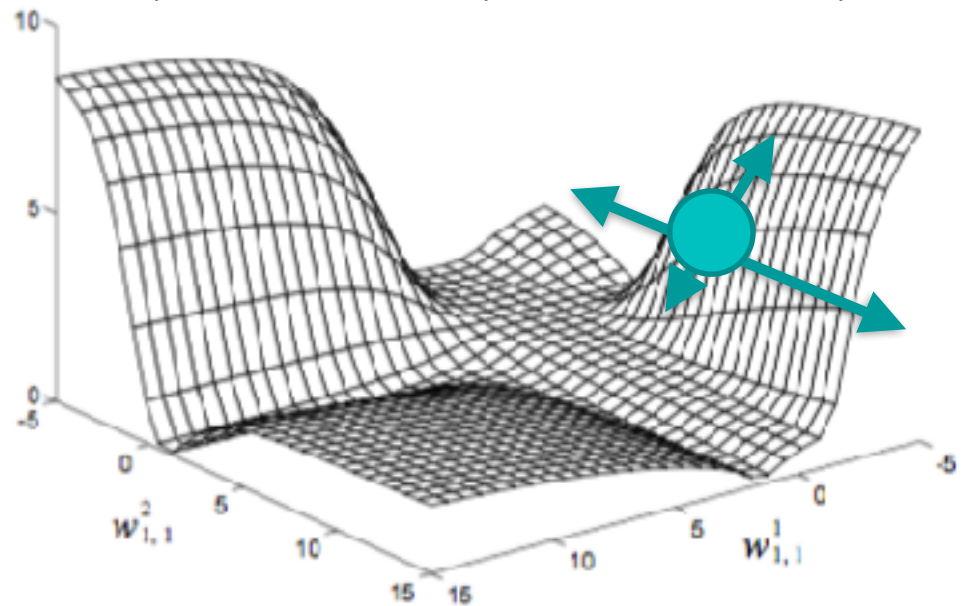
Be adaptive based on Gradient Magnitude?

Inverse magnitude of gradient in multiple directions?

$$\mathbf{W}_{k+1} \leftarrow \mathbf{W}_k + \eta \frac{1}{\sqrt{\mathbf{G}_k}} \odot \nabla J(\mathbf{W}_k)$$



$$\mathbf{G}_k = \nabla J(\mathbf{W}_k) \odot \nabla J(\mathbf{W}_k)$$



Common Adaptive Strategies $\mathbf{W}_{k+1} = \mathbf{W}_k - \eta \cdot \rho_k$

Adjust each element of gradient by the steepness

- AdaGrad
all operations are per element
$$\rho_k = \frac{1}{\sqrt{\mathbf{G}_k + \epsilon}} \odot \nabla J(\mathbf{W}_k) \quad \text{where} \quad \mathbf{G}_k = \mathbf{G}_{k-1} + \nabla J(\mathbf{W}_k) \odot \nabla J(\mathbf{W}_k)$$
- RMSProp
all operations are per element
$$\rho_k = \frac{1}{\sqrt{V_k + \epsilon}} \odot \nabla J(\mathbf{W}_k) \quad \begin{aligned} \mathbf{G}_k &= \nabla J(\mathbf{W}_k) \odot \nabla J(\mathbf{W}_k) \\ \mathbf{V}_k &= \gamma \cdot \mathbf{V}_{k-1} + (1 - \gamma) \cdot \mathbf{G}_k \end{aligned}$$
- AdaDelta
all operations are per element
$$\rho_k = \frac{\mathbf{M}_k}{\sqrt{\mathbf{V}_k + \epsilon}} \quad \mathbf{M}_{k+1} = \gamma \cdot \mathbf{M}_k + (1 - \gamma) \cdot \nabla J(\mathbf{W}_k)$$
- AdaM
 \mathbf{G} updates with decaying momentum of J and J^2
- NAdaM
same as Adam, but with nesterov's acceleration

None of these are “**one-size-fits-all**” because the space of neural network **optimization varies** by problem, ADAM is **popular** but **not a panacea**

Adaptive Momentum

All operations are element wise:

$$\beta_1 = 0.9, \beta_2 = 0.999, \eta = 0.001, \epsilon = 10^{-8}$$

$$k = 0, \mathbf{M}_0 = \mathbf{0}, \mathbf{V}_0 = \mathbf{0}$$

Published as a conference paper at ICLR 2015

ADAM: A METHOD FOR STOCHASTIC OPTIMIZATION

Diederik P. Kingma*
University of Amsterdam, OpenAI

Jimmy Lei Ba*
University of Toronto

For each epoch:

update epoch $k \leftarrow k + 1$

get gradient $\nabla J(\mathbf{W}_k)$

accumulated gradient $\mathbf{M}_k \leftarrow \beta_1 \cdot \mathbf{M}_{k-1} + (1 - \beta_1) \cdot \nabla J(\mathbf{W}_k)$

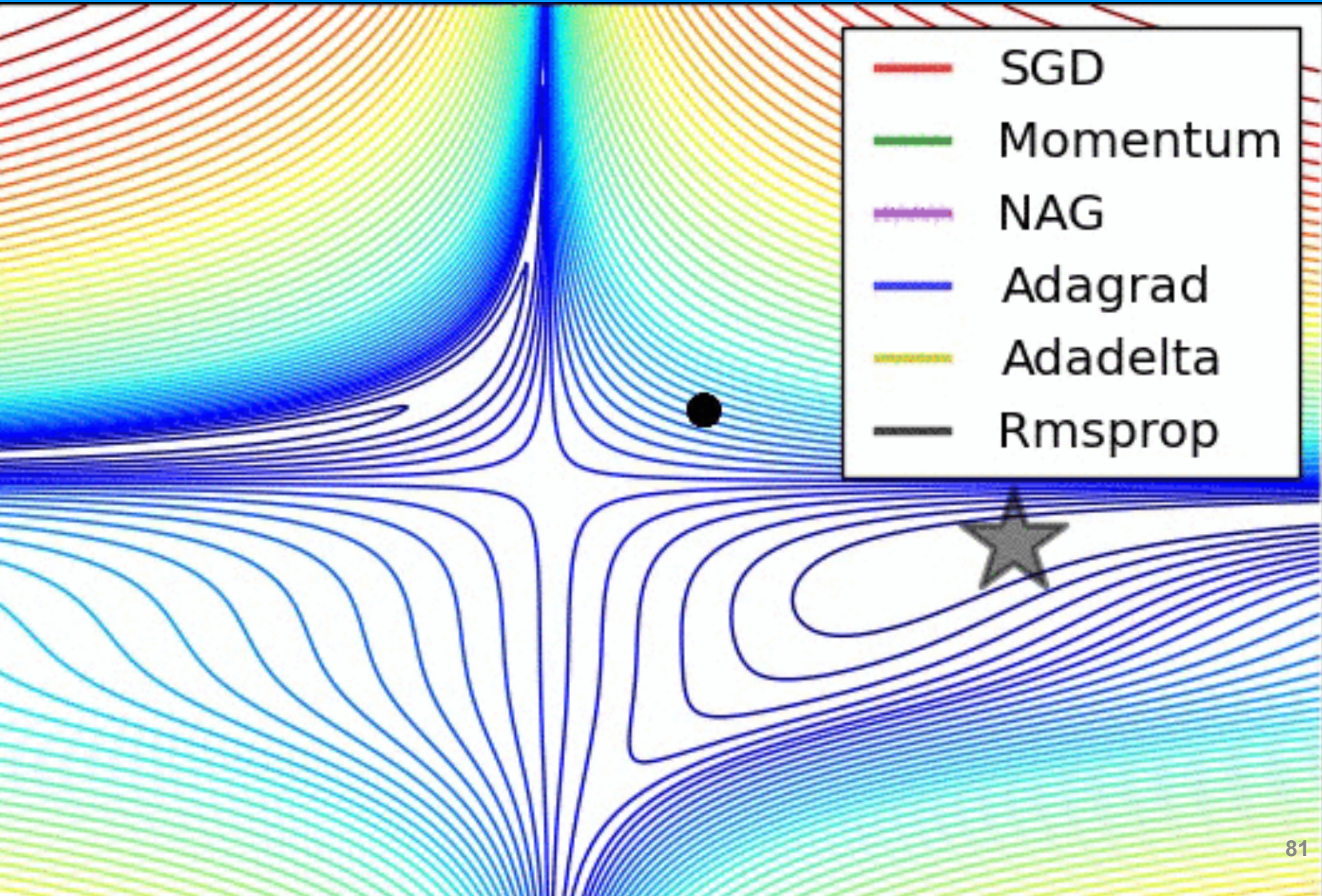
accumulated squared gradient $\mathbf{V}_k \leftarrow \beta_2 \cdot \mathbf{V}_{k-1} + (1 - \beta_2) \cdot \nabla J(\mathbf{W}_k) \odot \nabla J(\mathbf{W}_k)$

boost moments magnitudes
(notice k in exponent) $\hat{\mathbf{M}}_k \leftarrow \frac{\mathbf{M}_k}{(1 - [\beta_1]^k)} \quad \hat{\mathbf{V}}_k \leftarrow \frac{\mathbf{V}_k}{(1 - [\beta_2]^k)}$

update gradient, normalized
by second moment
similar to AdaDelta $\mathbf{W}_k \leftarrow \mathbf{W}_{k-1} - \eta \cdot \frac{\hat{\mathbf{M}}_k}{\sqrt{\hat{\mathbf{V}}_k + \epsilon}}$

Visualization of Optimization

<https://ruder.io/optimizing-gradient-descent/>



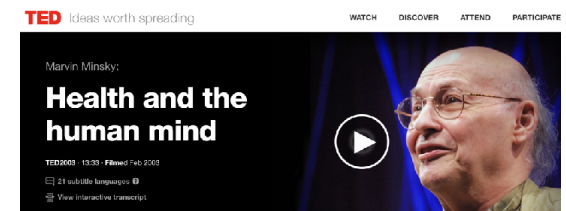
Practical Details

- Neural networks can separate any data through multiple layers. The true realization of Rosenblatt:

"Given an elementary α -perceptron, a stimulus world W , and any classification $C(W)$ for which a solution exists; let all stimuli in W occur in any sequence, provided that each stimulus must reoccur in finite time; then beginning from an arbitrary initial state, an error correction procedure will always yield a solution to $C(W)$ in finite time..."



- **Universality:** No matter what function we want to compute, we know that there is a neural network which can do the job.



- One nonlinear hidden layer with an output layer can perfectly train any problem with enough data, but might just be memorizing...
 - ... it might be better to have even more layers for decreased computation and generalizability

End of Session

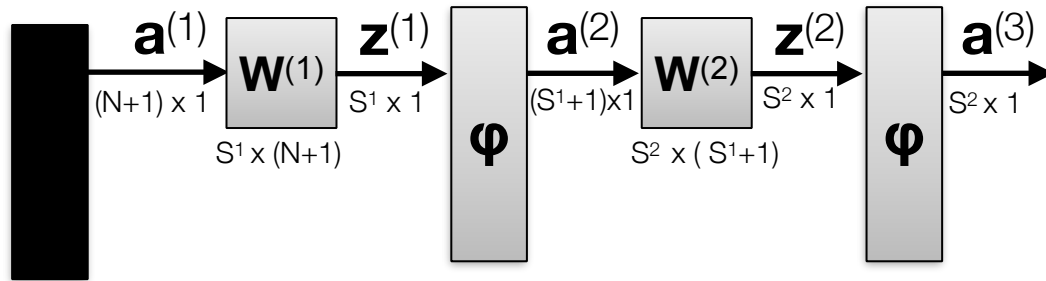
- Next Time: Final Flipped Module!
- Then: Deep Learning in Keras

Town Hall



Back Up Slides

Back propagation

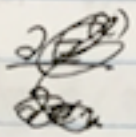
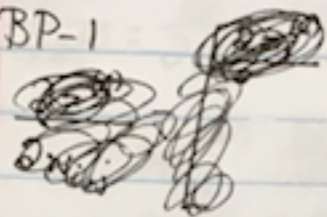


use chain rule:
$$\frac{\partial J(\mathbf{W})}{\partial w_{ij}^{(l)}} = \frac{\partial J(\mathbf{W})}{\partial \mathbf{z}^{(l)}} \frac{\partial \mathbf{z}^{(l)}}{\partial w_{ij}^{(l)}}$$

$$J(\mathbf{W}) = \sum_k^M (\mathbf{y}^{(k)} - \mathbf{a}^{(L)})^2$$

$$w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} - \eta \frac{\partial J(\mathbf{W})}{\partial w_{ij}^{(l)}}$$

BP-1



$$\frac{\partial Z^{(e)}}{\partial W_{i,j}^{(e)}} = \begin{bmatrix} \frac{\partial Z^{(e)}}{\partial W_{i,j}^{(e)}} \\ \vdots \\ \frac{\partial Z_{se}^{(e)}}{\partial W_{i,j}^{(e)}} \end{bmatrix}$$

MATRIX

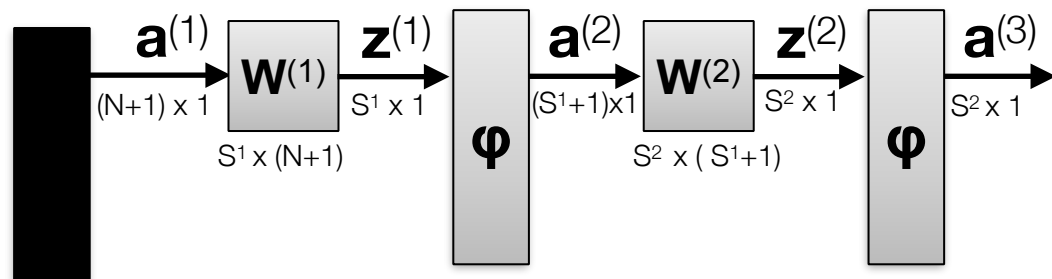
$$\begin{aligned} Z^{(e)} &= W^{(e)} a^{(e)} \\ Z_k^{(e)} &= \sum_{j=1}^n W_{kj}^{(e)} a_j^{(e)} \\ &= \sum_{j=1}^n W_{kj}^{(e)} a_j^{(e)} \end{aligned}$$

$$\begin{aligned} \frac{\partial Z_i^{(e)}}{\partial W_{i,j}^{(e)}} &= \sum_{j=1}^n \frac{\partial Z_i^{(e)}}{\partial W_{i,j}^{(e)}} a_j^{(e)} W_{k,j}^{(e)} \\ &= 0 \text{ if } k \neq i \end{aligned}$$

$$= a_j^{(e)} \text{ if } k = i$$

$$\frac{\partial Z^{(e)}}{\partial W_{i,j}^{(e)}} = a_j^{(e)} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ a_j^{(e)} \\ \vdots \\ 0 \end{bmatrix} \leftarrow i^{\text{TH}} \text{ POSITION}$$

Back propagation



can we use
chain rule
again?

$$\frac{\partial J(\mathbf{W})}{\partial w_{ij}^{(l)}} = \frac{\partial J(\mathbf{W})}{\partial \mathbf{z}^{(l)}} a_j^{(l)}$$

$$J(\mathbf{W}) = \sum_k^M (\mathbf{y}^{(k)} - \mathbf{a}^{(L)})^2$$

$$w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{z}^{(l)}} a_j^{(l)}$$

- Steps:
 - ❑ propagate weights forward
 - ❑ calculate gradient at final layer
 - ❑ back propagate gradient for each layer
 - ❑ via recurrence relation

BP-2, setup $J(w)/\partial z^{(l+1)}$

CHAIN RULE
AGAIN

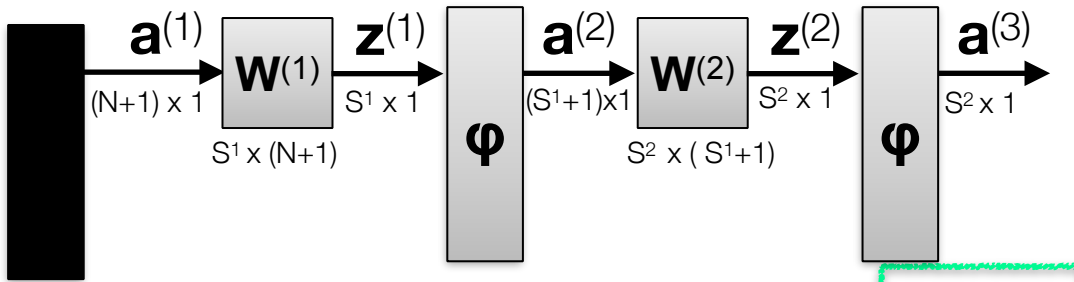
$$\frac{\partial J(w)}{\partial z^{(l)}} = \underbrace{\frac{\partial J(w)}{\partial z^{(l+1)}}}_{\text{VECTOR } 1 \times S^{(l+1)}} \underbrace{\frac{\partial z^{(l+1)}}{\partial z^{(l)}}}_{\text{MATRIX } S^{(l+1)} \times S^{(l)}}$$

$$\begin{bmatrix} \frac{\partial J(w)}{\partial z_1^{(l+1)}} \\ \vdots \\ \frac{\partial J(w)}{\partial z_{S^{(l+1)}}^{(l+1)}} \end{bmatrix} \begin{bmatrix} \frac{\partial z_1^{(l+1)}}{\partial z_1^{(l)}} & \frac{\partial z_2^{(l+1)}}{\partial z_1^{(l)}} & \dots & \frac{\partial z_{S^{(l+1)}}^{(l+1)}}{\partial z_1^{(l)}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial z_1^{(l+1)}}{\partial z_{S^{(l)}}^{(l)}} & \dots & \dots & \frac{\partial z_{S^{(l+1)}}^{(l+1)}}{\partial z_{S^{(l)}}^{(l)}} \end{bmatrix}$$

$$\frac{\partial z_i^{(l+1)}}{\partial z_j^{(l)}} = \frac{\partial}{\partial z_j^{(l)}} \sum_{R=1}^{S^{(l+1)}} w_{i,j,R}^{(l+1)} q_R^{(l+1)}$$

1/1

Back propagation



can we get this gradient?

$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{z}^{(l)}} = \frac{\partial J(\mathbf{W})}{\partial \mathbf{z}^{(l+1)}} \frac{\partial \mathbf{z}^{(l+1)}}{\partial \mathbf{z}^{(l)}}$$

$$J(\mathbf{W}) = \sum_k^M (\mathbf{y}^{(k)} - \mathbf{a}^{(L)})^2$$

$$w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{z}^{(l)}} a_j^{(l)}$$

BP-3

$$\frac{\partial z_i^{(e+1)}}{\partial z_j^{(e)}} = \frac{\partial}{\partial z_j^{(e)}} \left(\sum_{k=1}^{S^{e+1}} w_{i,k}^{(e+1)} q_k^{(e+1)} \right)$$

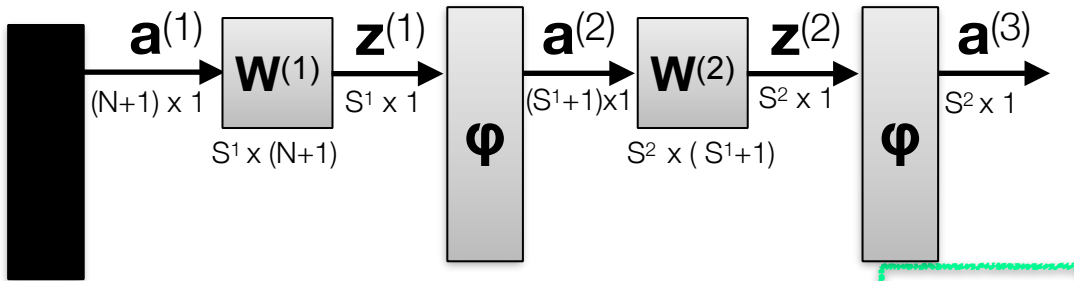
$$= w_{i,j}^{(e+1)} \frac{\partial}{\partial z_j^{(e)}} \phi(z_j^{(e)})$$

$$= w_{i,j}^{(e+1)} q_j^{(e+1)} (1 - q_j^{(e+1)})$$

$$\frac{\partial z^{(e+1)}}{\partial z^{(e)}} = \underbrace{W^{(e+1)} \text{diag}(q_j^{(e+1)} (1 - q_j^{(e+1)}))}_{S^{e+1} \times S^{e+1}} \underbrace{W^{(e+1)}}_{S^{e+1} \times S^e}$$

ENTIRE MATRIX
 $S^{e+1} \times S^e$

Back propagation



use chain rule:
$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{z}^{(l)}} = \frac{\partial J(\mathbf{W})}{\partial \mathbf{z}^{(l+1)}} \frac{\partial \mathbf{z}^{(l+1)}}{\partial \mathbf{z}^{(l)}}$$

$$\frac{\partial \mathbf{z}^{(l+1)}}{\partial \mathbf{z}^{(l)}} = \text{diag}[\mathbf{a}^{(l+1)} * (1 - \mathbf{a}^{(l+1)})] \cdot \mathbf{W}^{(l+1)}$$

$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{z}^{(l)}} = \text{diag}[\mathbf{a}^{(l+1)} * (1 - \mathbf{a}^{(l+1)})] \cdot \mathbf{W}^{(l+1)} \frac{\partial J(\mathbf{W})}{\partial \mathbf{z}^{(l+1)}}$$

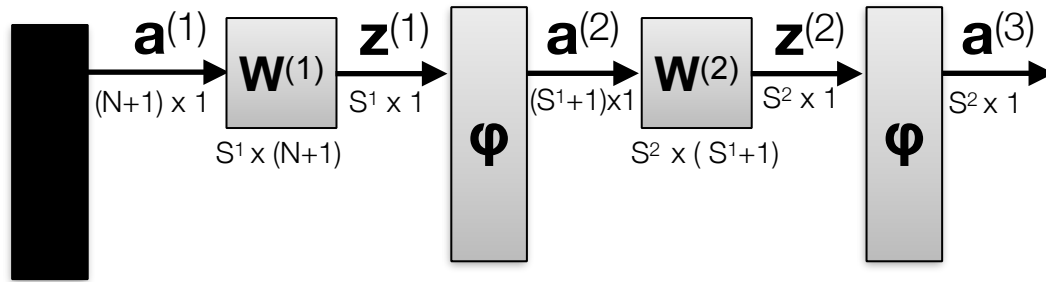
recurrence relation

$$J(\mathbf{W}) = \sum_k^M (\mathbf{y}^{(k)} - \mathbf{a}^{(L)})^2$$

$$w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{z}^{(l)}} a_j^{(l)}$$

If we know last layer, we can **back propagate** towards previous layers!

Back propagation



one more step
need last layer
gradient:

$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{z}^{(l+1)}}$$

$$J(\mathbf{W}) = \sum_k^M (\mathbf{y}^{(k)} - \mathbf{a}^{(L)})^2$$

$$w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{z}^{(l)}} a_j^{(l)}$$

$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{z}^{(2)}} = \frac{\partial}{\partial \mathbf{z}^{(2)}} (\mathbf{y}^{(k)} - \phi(\mathbf{z}^{(2)}))^2$$

$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{z}^{(2)}} = -2(\mathbf{y}^{(k)} - \mathbf{a}^{(3)}) * \mathbf{a}^{(3)} * (1 - \mathbf{a}^{(3)})$$

Final Layer BP-4

$$\frac{\partial J(W)}{\partial z_i^{(L)}} = \frac{\partial}{\partial z_i^{(L)}} \sum_{k=1}^M (y^{(k)} - \phi(z^{(L)}))^2$$

$$= - \sum_k 2(y^{(k)} - a^{(L+1)}) \frac{\partial}{\partial z_i^{(L)}} \phi(z^{(L)})$$

$$= -2 \sum_k (y^{(k)} - a^{(L+1)}) a_i^{(L+1)} (1 - a_i^{(L+1)}) \frac{\partial}{\partial z_i^{(L)}} z^{(L)} \quad \textcircled{1}$$

$$= -2 \sum_k (y^{(k)} - a^{(L+1)}) a_i^{(L+1)} (1 - a_i^{(L+1)})$$

$$\frac{\partial J}{\partial z^{(L)}} = -2 \sum_k (y^{(k)} - a^{(L+1)}) a^{(L+1)} (1 - a^{(L+1)})$$

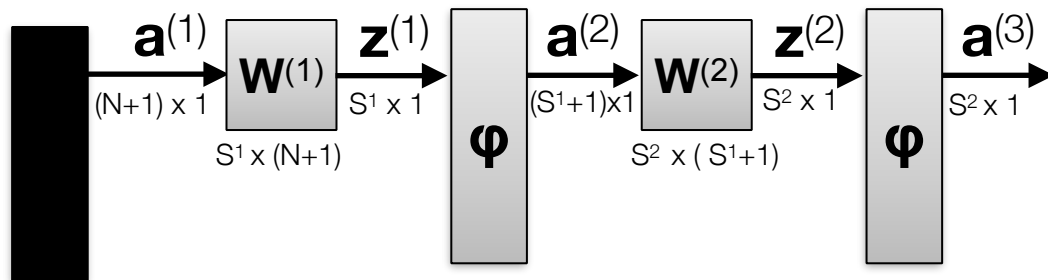
↑ ONE DIM, TAKE OUT K

OR w/ layers

$$\frac{\partial J(W)}{\partial z^{(2)}} = \frac{\partial}{\partial z^{(2)}} \sum_k (y^{(k)} - \phi(z^{(2)}))^2$$

$$= -2 \sum_k (y^{(k)} - a^{(3)}) a^{(3)} (1 - a^{(3)})$$

Back propagation summary



- **Self Test:**

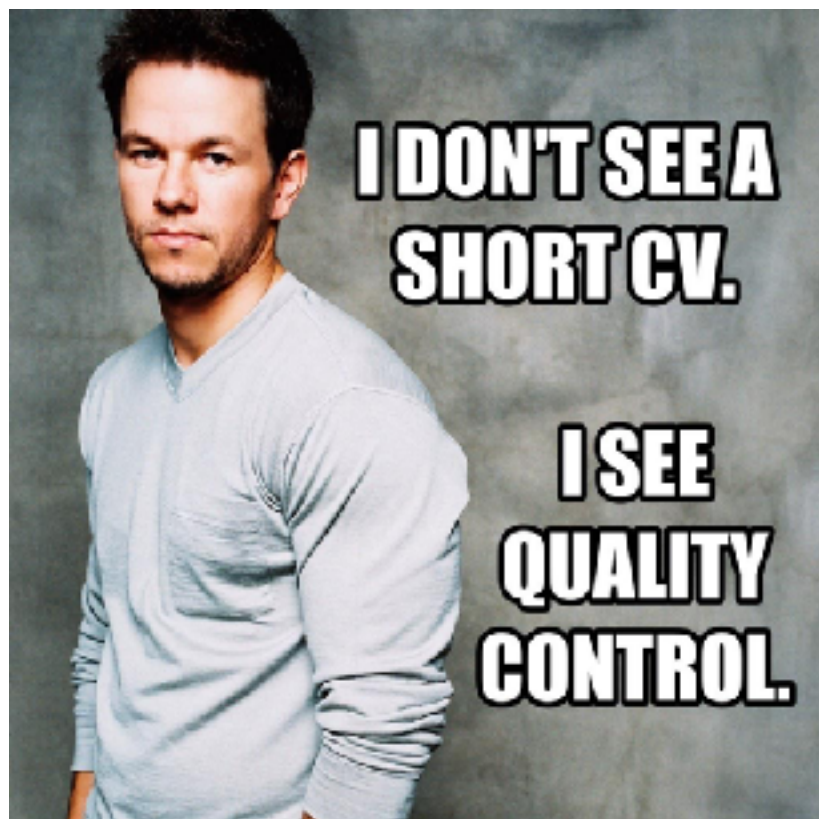
True or False: If we change the cost function, $J(\mathbf{W})$, we only need to update the final layer calculation of the back propagation steps. The remainder of the algorithm is unchanged.

- A. True
- B. False

$$\mathbf{V}^{(2)} = -2(\mathbf{Y} - \mathbf{A}^{(3)}) * \mathbf{A}^{(3)} * (1 - \mathbf{A}^{(3)})$$
$$\nabla^{(2)} = \mathbf{V}^{(2)} \cdot [\mathbf{A}^{(2)}]^T$$

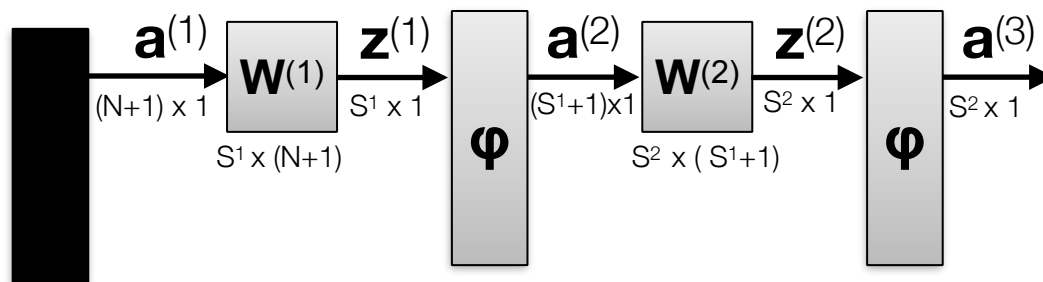
Programming Multi-layer Neural Networks

07. MLP Neural Networks.ipynb



Guided Example

Recall from the in-class assignment



1. Forward propagate to get \mathbf{Z} , \mathbf{A}

```
# feedforward all instances
```

```
A1, Z1, A2, Z2, A3 = self._feedforward(X_data, self.W1, self.W2)
```

```
def _feedforward(self, X, W1, W2):
```

```
    A1 = self._add_bias_unit(X.T, how='row')
```

```
    Z1 = W1 @ A1
```

```
    A2 = self._sigmoid(Z1)
```

```
    A2 = self._add_bias_unit(A2, how='row')
```

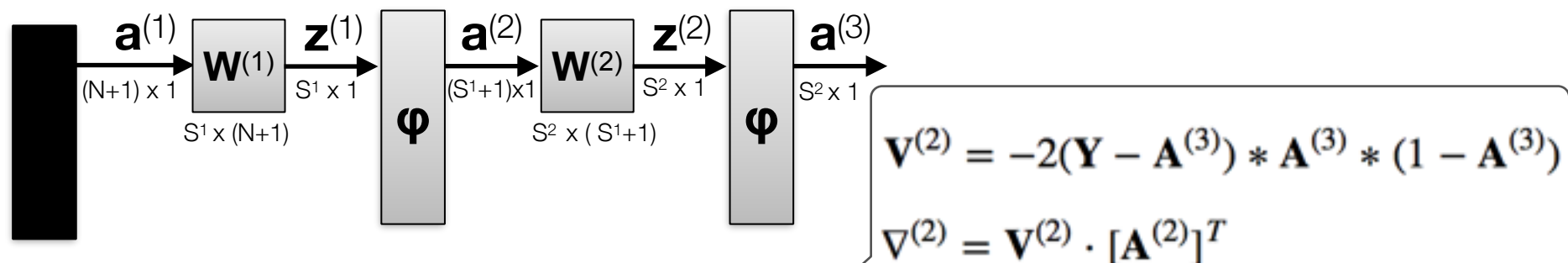
```
    Z2 = W2 @ A2
```

```
    A3 = self._sigmoid(Z2)
```

```
    return A1, Z1, A2, Z2, A3
```

these are more than just vectors for **one instance**!
these are for **all** instances

Back propagation implementation



1. Forward propagate to get \mathbf{Z} , \mathbf{A}
2. Get final layer gradient
3. Back propagate sensitivities

$$\mathbf{V}^{(2)} = -2(\mathbf{Y} - \mathbf{A}^{(3)}) * \mathbf{A}^{(3)} * (1 - \mathbf{A}^{(3)})$$

$$\nabla^{(2)} = \mathbf{V}^{(2)} \cdot [\mathbf{A}^{(2)}]^T$$

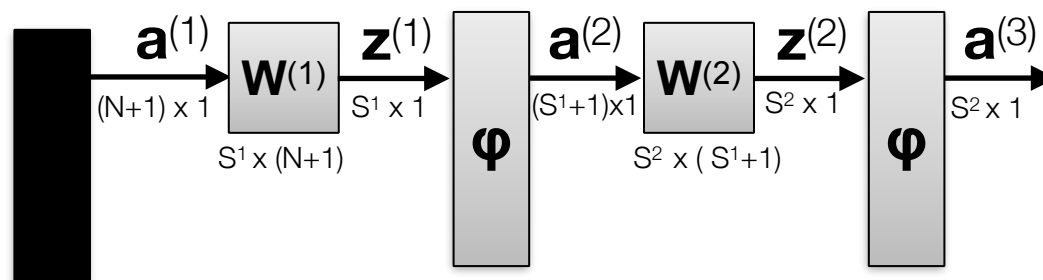
$$\mathbf{V}^{(1)} = \mathbf{A}^{(2)} * (1 - \mathbf{A}^{(2)}) * [\mathbf{W}^{(2)}]^T \cdot \mathbf{V}^{(2)}$$

$$\nabla^{(1)} = \mathbf{V}^{(1)} \cdot [\mathbf{A}^{(1)}]^T$$

```
def _get_gradient(self, A1, A2, A3, Z1, Z2, Y_enc, W1, W2):
    """ Compute gradient step using backpropagation.
    """
    # vectorized backpropagation
    V2 = -2*(Y_enc-A3)*A3*(1-A3)
    V1 = A2*(1-A2)*(W2.T @ V2)

    grad2 = V2 @ A2.T
    grad1 = V1[1:,:] @ A1.T
```

Back propagation implementation



1. Forward propagate to get \mathbf{Z} , \mathbf{A} for all layers
2. Get final layer gradient
3. Back propagate sensitivities
4. Update each $\mathbf{W}^{(l)}$

for each layer:

$$\mathbf{W}^{(l)} \leftarrow \mathbf{W}^{(l)} - \eta \nabla^{(l)}$$

```
# feedforward all instances
A1, Z1, A2, Z2, A3 = self._feedforward(X_data, self.W1, self.W2)

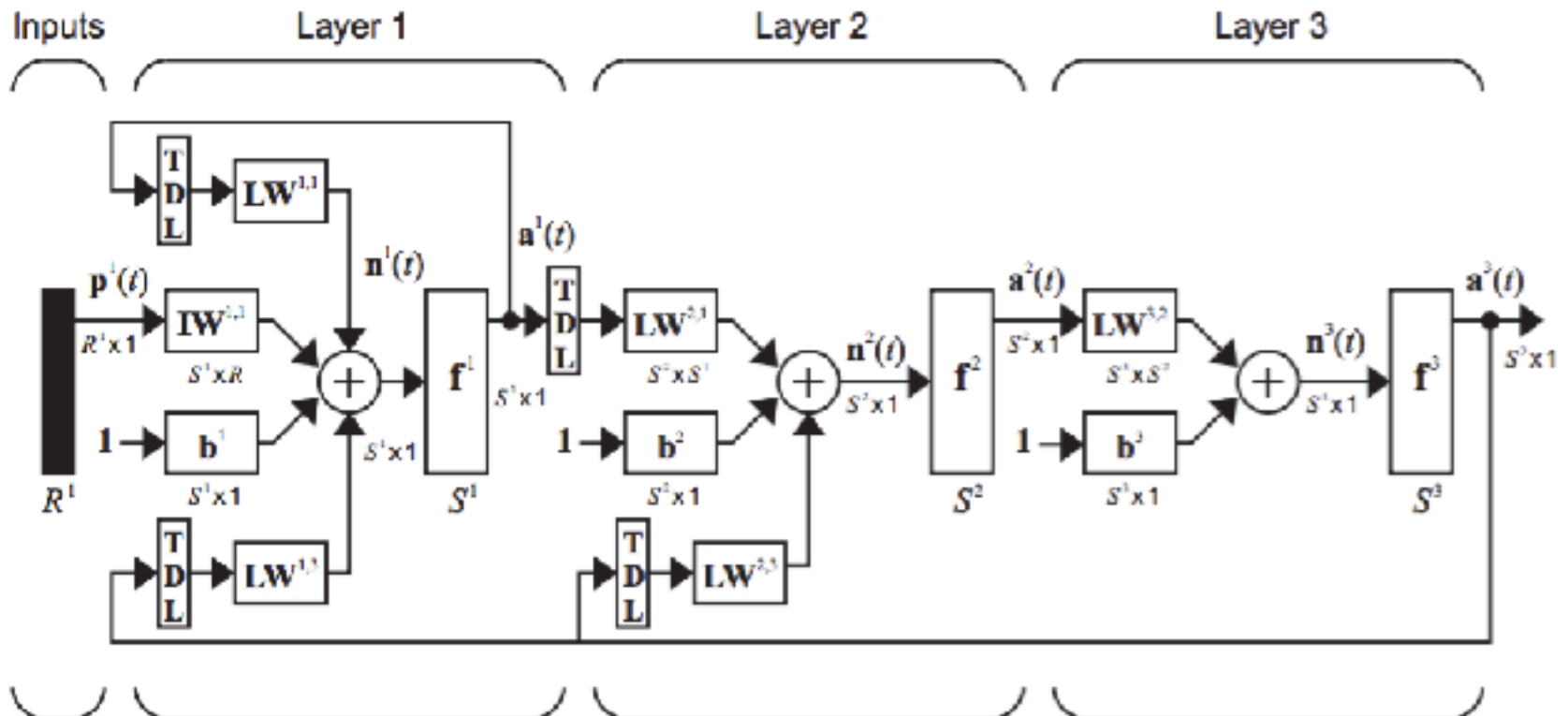
# compute gradient via backpropagation
grad1, grad2 = self._get_gradient(A1=A1, A2=A2,
                                   A3=A3, Z1=Z1,
                                   Y_enc=Y_enc,
                                   W1=self.W1, W2=self.W2)

self.W1 -= self.eta * grad1
self.W2 -= self.eta * grad2
```

More Advanced Architectures

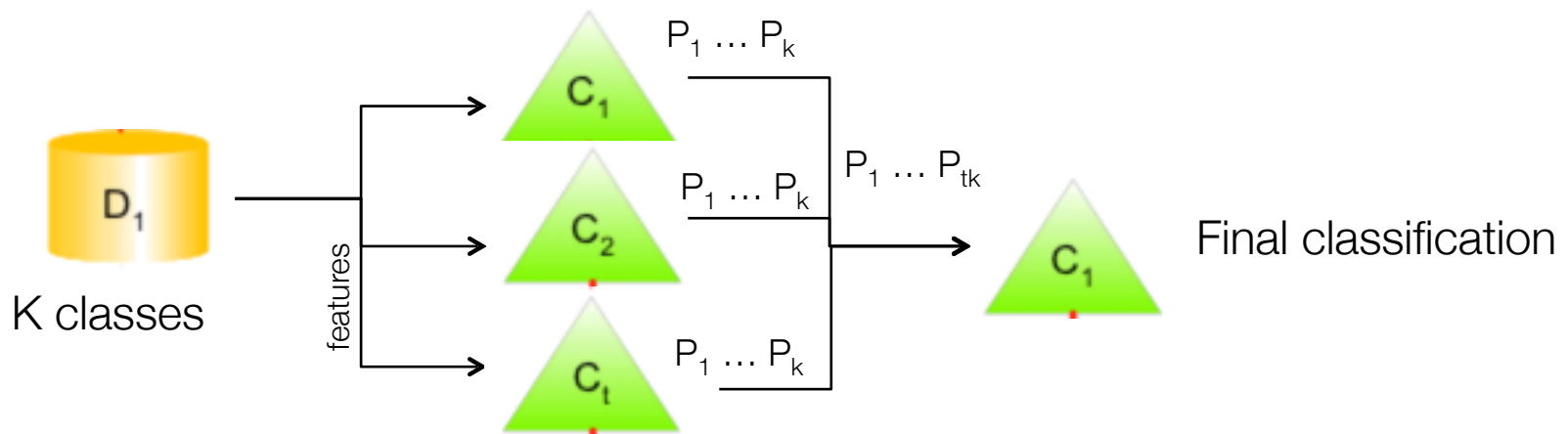
- Dynamic Networks (recurrent networks)
 - can use current and previous inputs, in time
 - still popular, but ultimately extremely hard to train
 - highly successful variant:** long short term

1. LSTM



Stacking and Cascading ensembles

- Train an initial classifier (or ensemble)
- Use the output probabilities of each class (from each classifier, if ensemble), as features for another classifier
- Train final classifier on inferences from previous classes



Self test

- How many features (attributes) does scikit-learn use in each iteration of bagging for a random forest?
- Which best completes this statement: In scikit, the number of features considered changes
 - ☐ a: once for each tree in the forest
 - ☐ b: once for each split in each tree

each iteration!!!

```
max_features=sqrt(n_features)
```


SVMs Summary: Linear

- Linear SVMs
 - ▣ Architecture identical to logistic regression, except:
 - ▣ maximize margin
 - ▣ constrained optimization
 - ▣ **Self Test:** What are the trained parameters for the linear SVM?

▣ A: Coefficients of w

▣ B: Intercept

▣ C: Slack Variables

▣ D: Support Vector locations

$$\begin{aligned} \min_{w, b, \zeta} & \frac{1}{2} w^T w + C \sum_{i=1}^n \zeta_i \\ \text{to } & y_i(w^T \phi(x_i) + b) \geq 1 - \zeta_i, \\ & \zeta_i \geq 0, i = 1, \dots, n \end{aligned}$$

SVMs Summary: Linear

- Linear SVMs

- ▢ Architecture near identical to logistic regression, except:

- ▢ maximize margin

- ▢ constrained optimization

$$\begin{aligned} \min_{w,b,\zeta} \quad & \frac{1}{2} w^T w + C \sum_{i=1}^n \zeta_i \\ \text{to} \quad & y_i (w^T \phi(x_i) + b) \geq 1 - \zeta_i, \\ & \zeta_i \geq 0, i = 1, \dots, n \end{aligned}$$

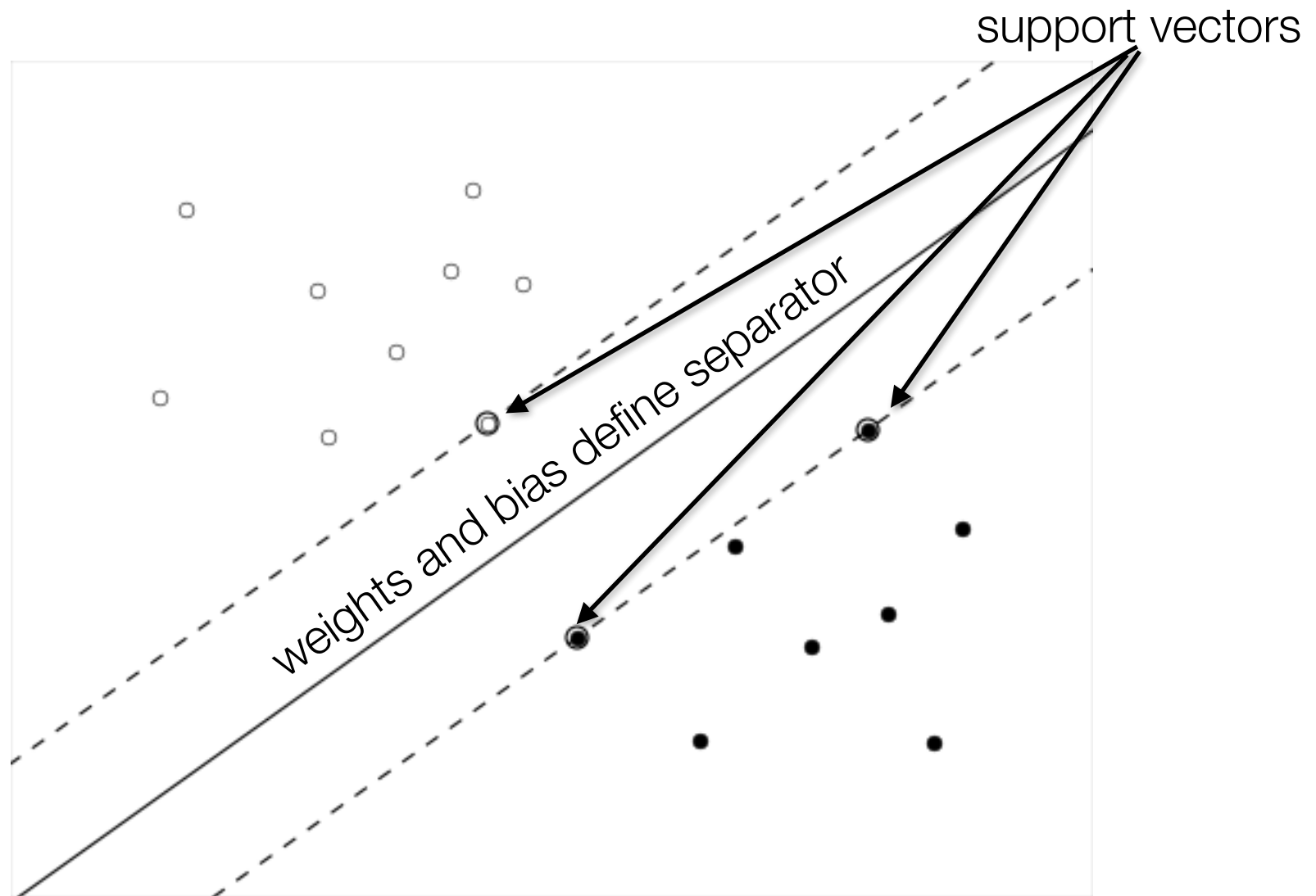
- ▢ Trained Parameters:

- ▢ intercept and weights for each class

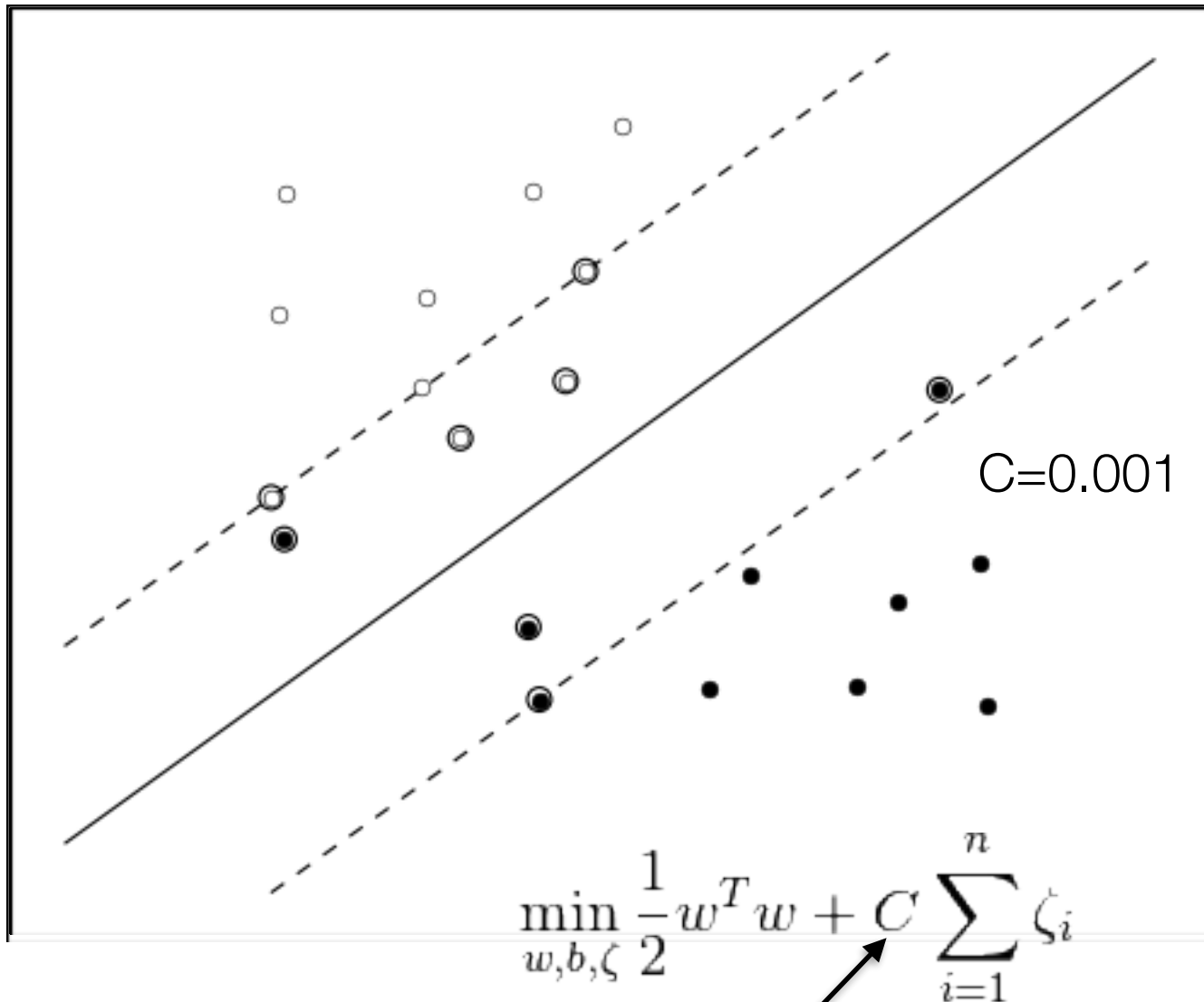
- ▢ support vectors chosen for margin calculation

- ▢ slack variables only needed during training

SVMs Summary: Linear

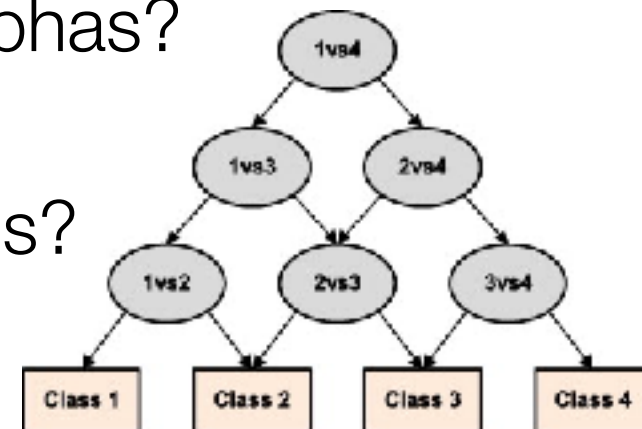


SVMs Summary: Linear



SVMs Summary: non-linear

- ❑ Non-linear SVMs
 - ❑ Architecture not like logistic regression
 - ❑ kernels == high dimensional dot-product
 - ❑ impossible to store weights
 - ❑ use kernel trick so no need to store them!
 - ❑ Trained Parameters
 - ❑ biases? how many?
 - ❑ selected support vectors? alphas?
 - ❑ slack variables?
 - ❑ parameters specific to kernels?



SVMs Summary: non-linear

Popular Kernels

polynomial

$$(\gamma \langle x, x' \rangle + r)^d$$

Diagram illustrating the polynomial kernel formula: $(\gamma \langle x, x' \rangle + r)^d$. Arrows point from the labels "gamma", "coef0", and "degree" to the corresponding terms in the formula: γ (gamma), r (coef0), and d (degree).

radial basis function

$$\exp(-\gamma |x - x'|^2)$$

Diagram illustrating the radial basis function kernel formula: $\exp(-\gamma |x - x'|^2)$. An arrow points from the label "gamma" to the γ term in the formula.

sigmoid

$$\tanh(\gamma \langle x, x' \rangle + r)$$

Diagram illustrating the sigmoid kernel formula: $\tanh(\gamma \langle x, x' \rangle + r)$. Arrows point from the labels "gamma" and "coef0" to the corresponding terms in the formula: γ (gamma) and r (coef0).

svm_gui.py



The End of SVMs