

# Distributed Supervised Machine Learning

Stony Brook University  
CSE545, Spring 2019

# Supervised Learning

(genes)

$X_1$

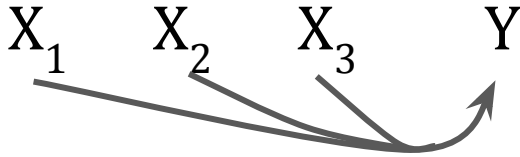
$X_2$

$X_3$

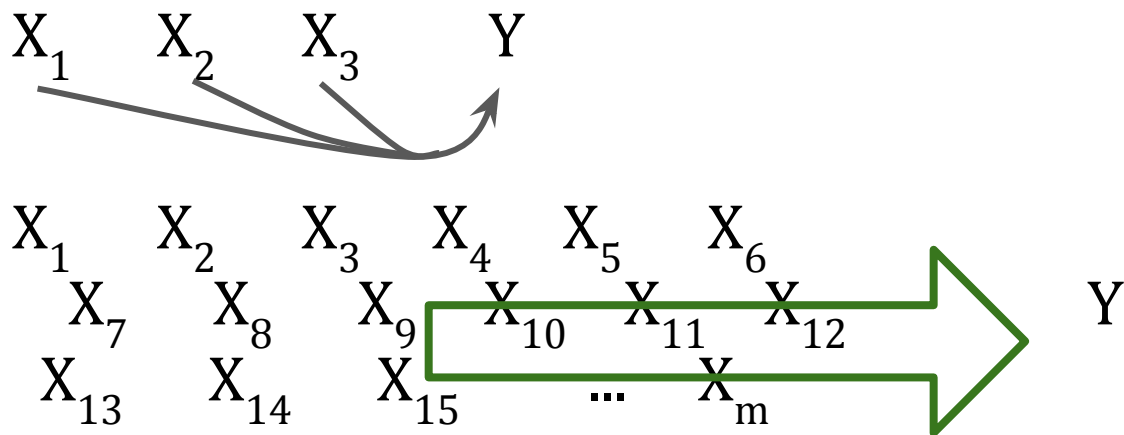
(health)

$Y$

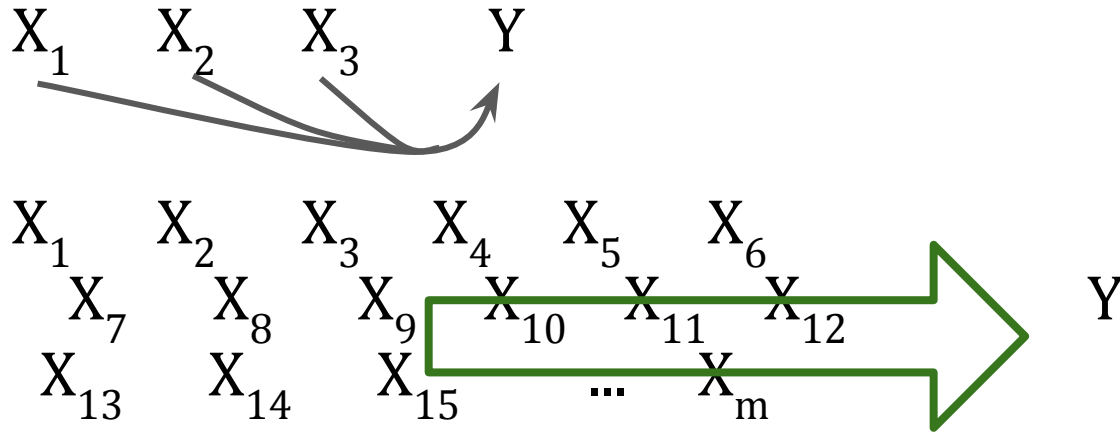
# Supervised Learning



# Supervised Learning

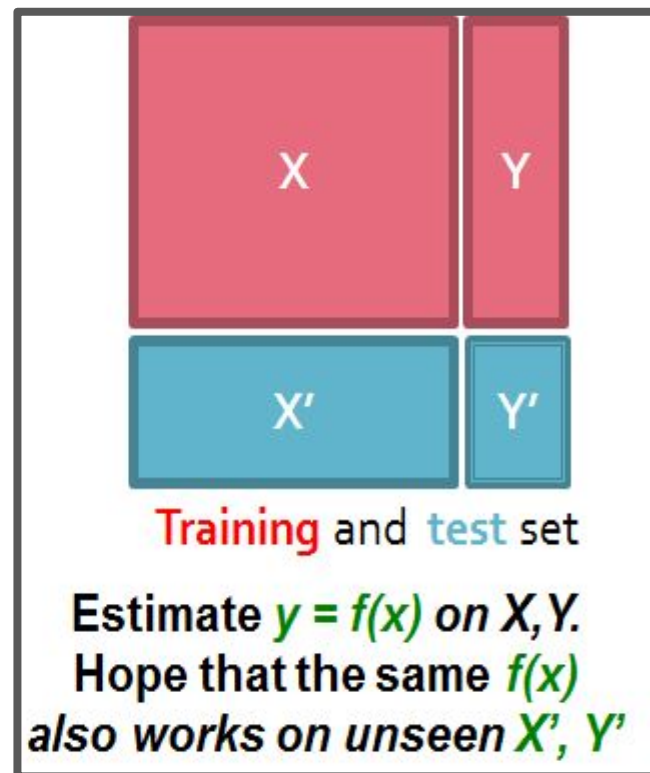
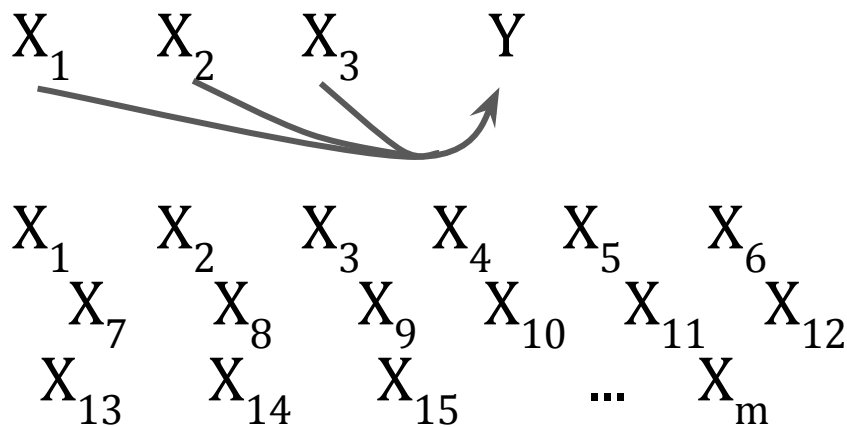


# Supervised Learning



Task: Determine a function,  $f$  (or parameters to a function) such that  $f(X) = Y$

# Supervised Learning



J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, <http://www.mmms.org>

Task: Determine a function,  $f$  (or parameters to a function) such that  $f(X) = Y$

# Supervised Learning Approaches that we will cover

1. **Regularized linear modeling**  
(linear and logistic regression)
2. Convolutional Neural Networks  
Where  $X$  might have spatial relationships
3. Recurrent Neural Networks  
Where  $X$  is a sequence of data

# Linear Model -- Linear Regression

Finding a linear function based on  $X$  to best yield  $Y$ .

$X$  = “covariate” = “feature” = “predictor” = “regressor” = “independent variable”

$Y$  = “response variable” = “outcome” = “dependent variable”

Regression:  $r(x) = E(Y|X = x)$

goal: estimate the function  $r$



$$r(x) = \beta_0 + \beta_1 x$$

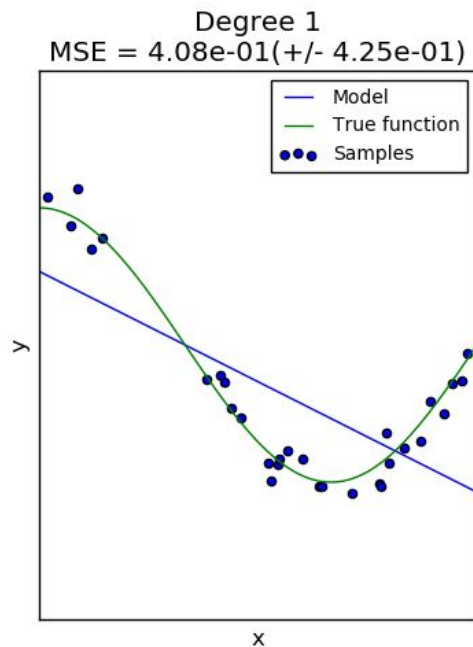
## Uses of linear and logistic regression $r(x) \approx E(Y|X = x)$

1. Testing the relationship between variables given other variables.  $\beta$  is an “effect size” -- a score for the magnitude of the relationship; can be tested for significance.
2. Building a predictive model that generalizes to new data.  $\hat{Y}$  is an estimate value of  $Y$  given  $X$ .

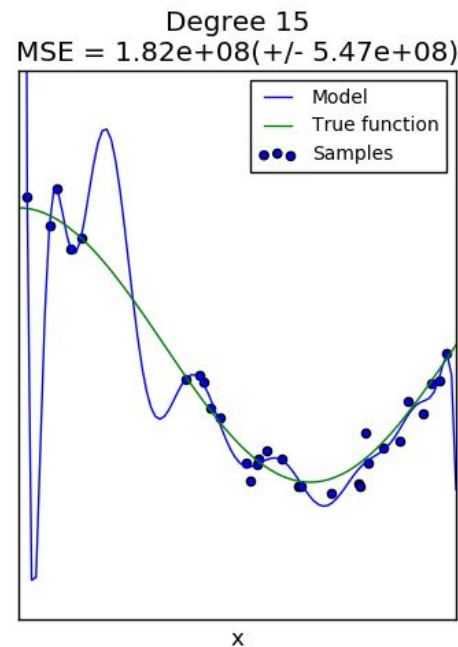
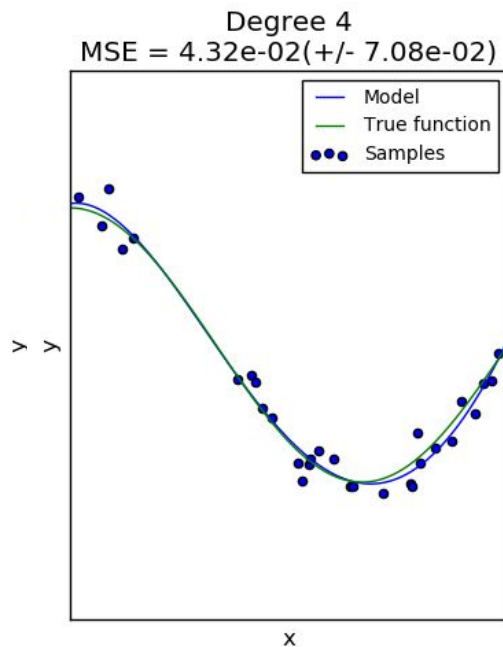
# Uses of linear and logistic regression

1. Testing the relationship between variables given other variables.  $\beta$  is an “effect size” -- a score for the magnitude of the relationship; can be tested for significance.
2. Building a predictive model that generalizes to new data.  $\hat{Y}$  is an estimate value of  $Y$  given  $X$ .  
However, unless  $|X| \ll \text{observations}$  then the model might “overfit”.

# Overfitting (1-d non-linear example)



Underfit  
High Bias



Overfit  
High Variance

*(image credit: Scikit-learn; in practice data are rarely this clear)*

# Overfitting (6-d linear example)

$$Y = X$$

1	0.5	0	0.6	1	0	0.25
1	0	0.5	0.3	0	0	0
0	0	0	1	1	1	0.5
0	0	0	0	0	1	1
1	0.25	1	1.25	1	0.1	2

# Overfitting (5-d linear example)

$$Y = X$$

1	0.5	0	0.6	1	0	0.25
1	0	0.5	0.3	0	0	0
0	0	0	1	1	1	0.5
0	0	0	0	0	1	1
1	0.25	1	1.25	1	0.1	2

$$\text{logit}(Y) = 1.2 + -63*X_1 + 179*X_2 + 71*X_3 + 18*X_4 + -59*X_5 + 19*X_6$$

# Overfitting (5-d linear example)

Do we really think we found something generalizable?

$Y$	=	$X$					
1		0.5	0	0.6	1	0	0.25
1		0	0.5	0.3	0	0	0
0		0	0	1	1	1	0.5
0		0	0	0	0	1	1
1		0.25	1	1.25	1	0.1	2

$$\text{logit}(Y) = 1.2 + -63*X_1 + 179*X_2 + 71*X_3 + 18*X_4 + -59*X_5 + 19*X_6$$

## Overfitting (2-d linear example)

Do we really think we found something generalizable?

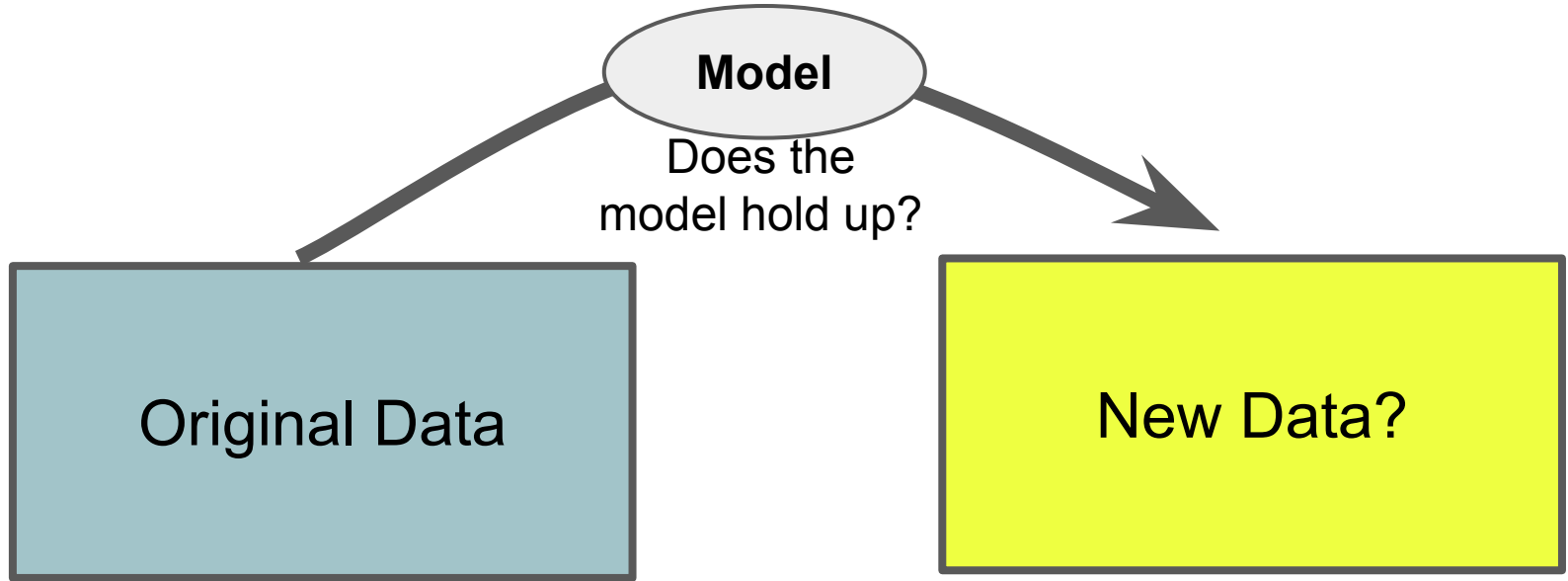
$$Y = X$$

1	0.5	0
1	0	0.5
0	0	0
0	0	0
1	0.25	1

What if only 2 predictors?

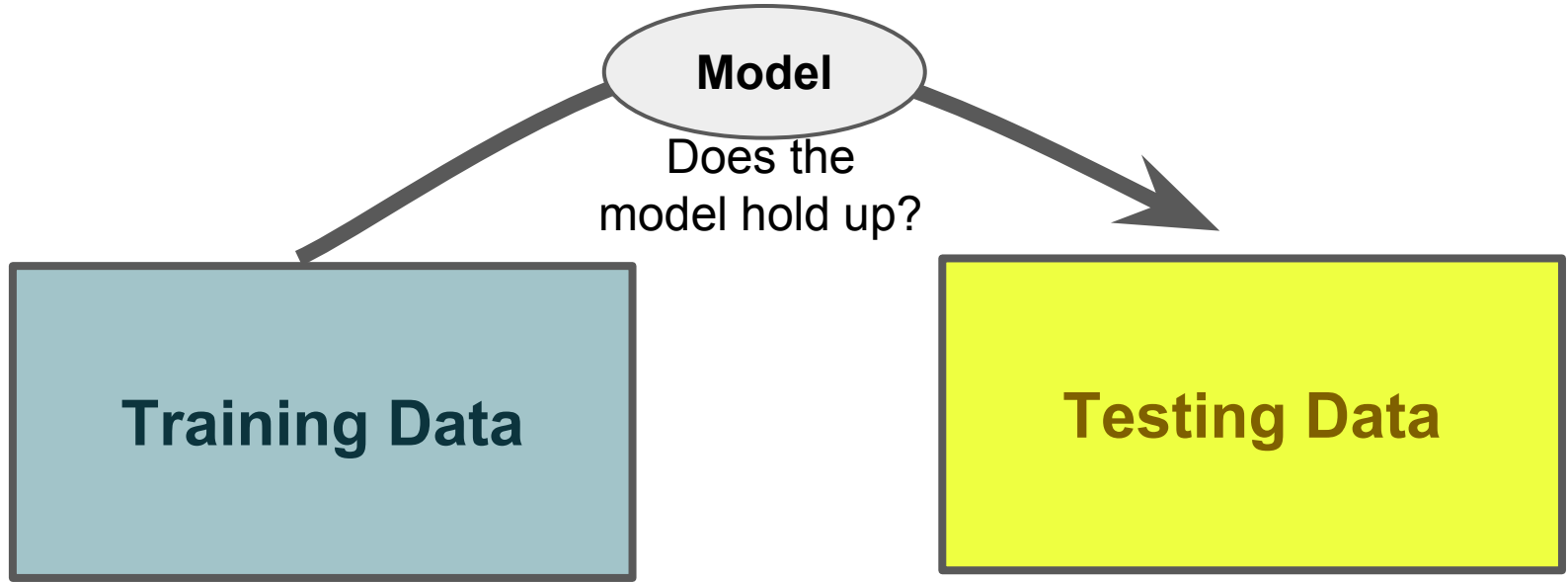
$$\text{logit}(Y) = 0 + 2 * X_1 + 2 * X_2$$

# Common Goal: Generalize to new data





# Common Goal: Generalize to new data





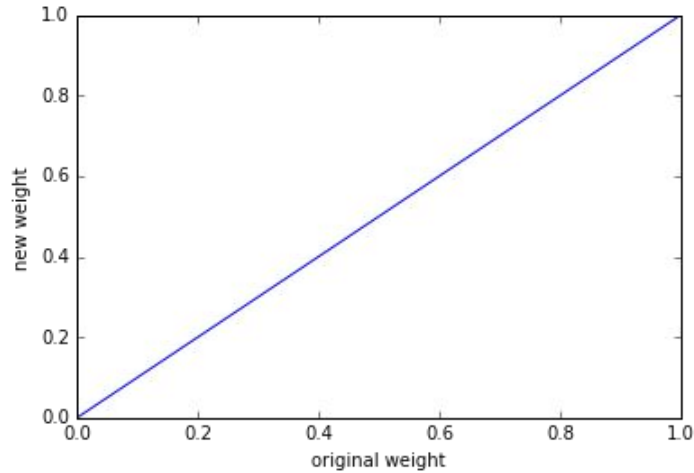
# Feature Selection / Subset Selection

## (bad) solution to overfit problem

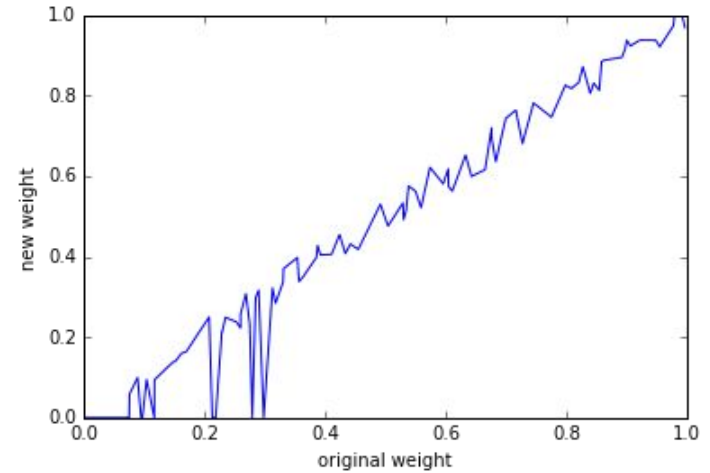
Use less features based on Forward Stepwise Selection:

- start with `current_model` just has the intercept (mean)  
`remaining_predictors = all_predictors`  
for `i` in `range(k)`:  
    #find best `p` to add to `current_model`:  
    for `p` in `remaining_predictors`  
        refit `current_model` with `p`  
        #add best `p`, based on  $RSS_p$  to `current_model`  
    #remove `p` from `remaining predictors`

# Regularization (Shrinkage)



No selection (weight= $\beta$ )



forward stepwise

Why just keep or discard features?

# Regularization (L2, Ridge Regression)

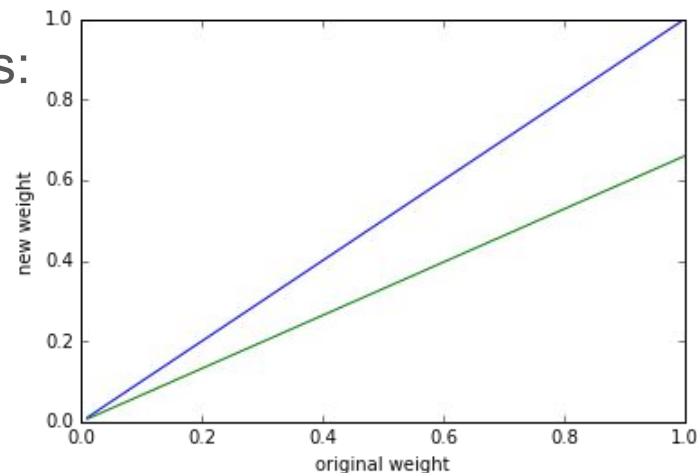
Idea: Impose a penalty on size of weights:

Ordinary least squares objective:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N (y_i - \sum_{j=1}^m x_{ij} \beta_j)^2 \right\}$$

Ridge regression:

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N (y_i - \sum_{j=1}^m x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^m \beta_j^2 \right\}$$



# Regularization (L2, Ridge Regression)

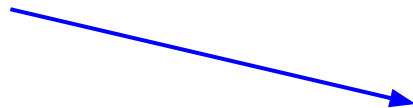
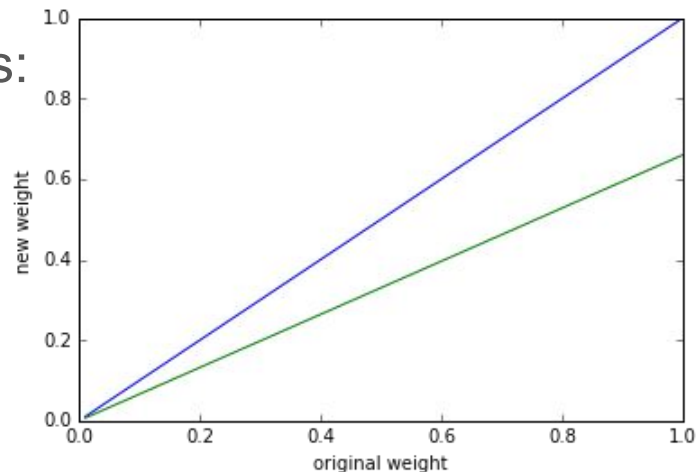
Idea: Impose a penalty on size of weights:

Ordinary least squares objective:

$$\hat{\beta} = \operatorname{argmin}_{\beta} \left\{ \sum_{i=1}^N (y_i - \sum_{j=1}^m x_{ij} \beta_j)^2 \right\}$$

Ridge regression:

$$\hat{\beta}^{\text{ridge}} = \operatorname{argmin}_{\beta} \left\{ \sum_{i=1}^N (y_i - \sum_{j=1}^m x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^m \beta_j^2 \right\}$$



$$\lambda \|\beta\|_2^2$$

# Regularization (L2, Ridge Regression)

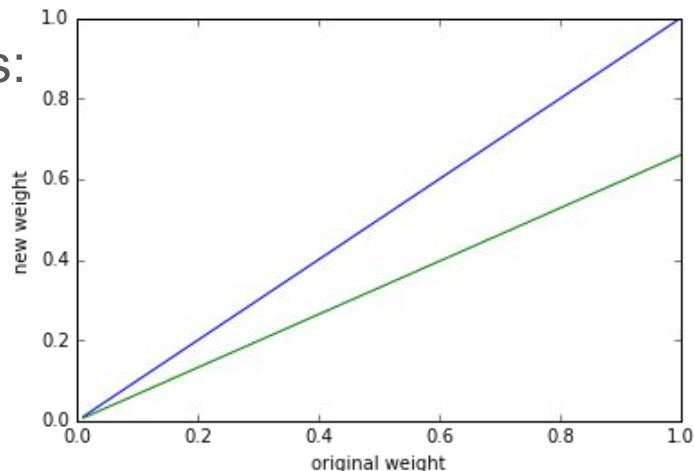
Idea: Impose a penalty on size of weights:

Ordinary least squares objective:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N (y_i - \sum_{j=1}^m x_{ij} \beta_j)^2 \right\}$$

Ridge regression:

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N (y_i - \sum_{j=1}^m x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^m \beta_j^2 \right\}$$



In Matrix Form:

$$\text{RSS}(\lambda) = (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta$$

$$\hat{\beta}^{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T y$$

$I$ :  $m \times m$  identity matrix

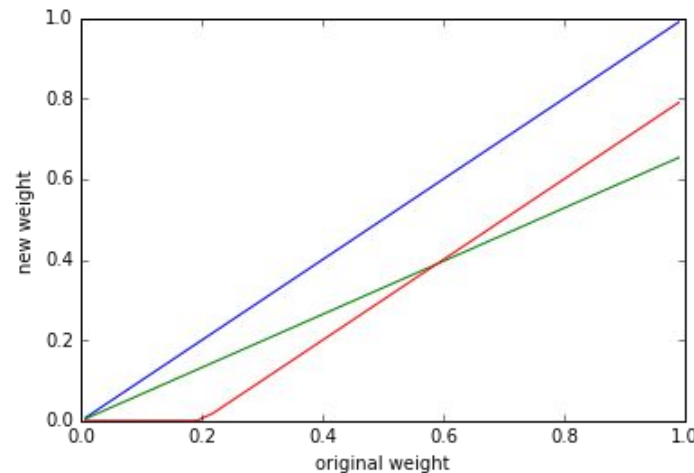
$$\lambda \|\beta\|_2^2$$

# Regularization (L1, The “Lasso”)

Idea: Impose a penalty and zero-out some weights

The Lasso Objective:

$$\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2} \sum_{i=1}^N (Y_i - \sum_{j=1}^m x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^m |\beta_j| \right\}$$



$\lambda \|\beta\|_1$

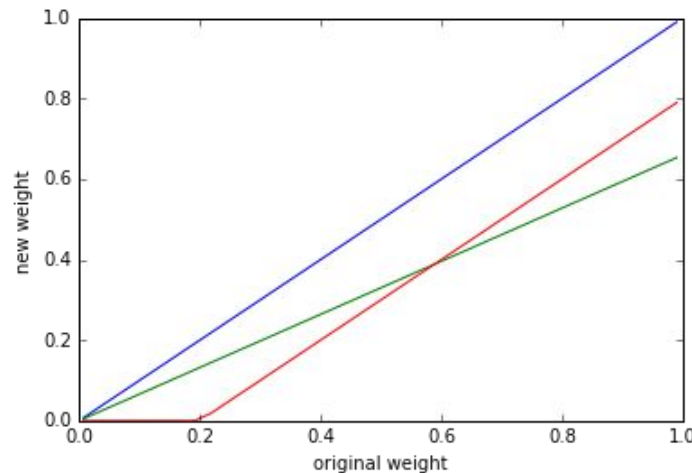


# Regularization (L1, The “Lasso”)

Idea: Impose a penalty and zero-out some weights

The Lasso Objective:

$$\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2} \sum_{i=1}^N (Y_i - \sum_{j=1}^m x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^m |\beta_j| \right\}$$

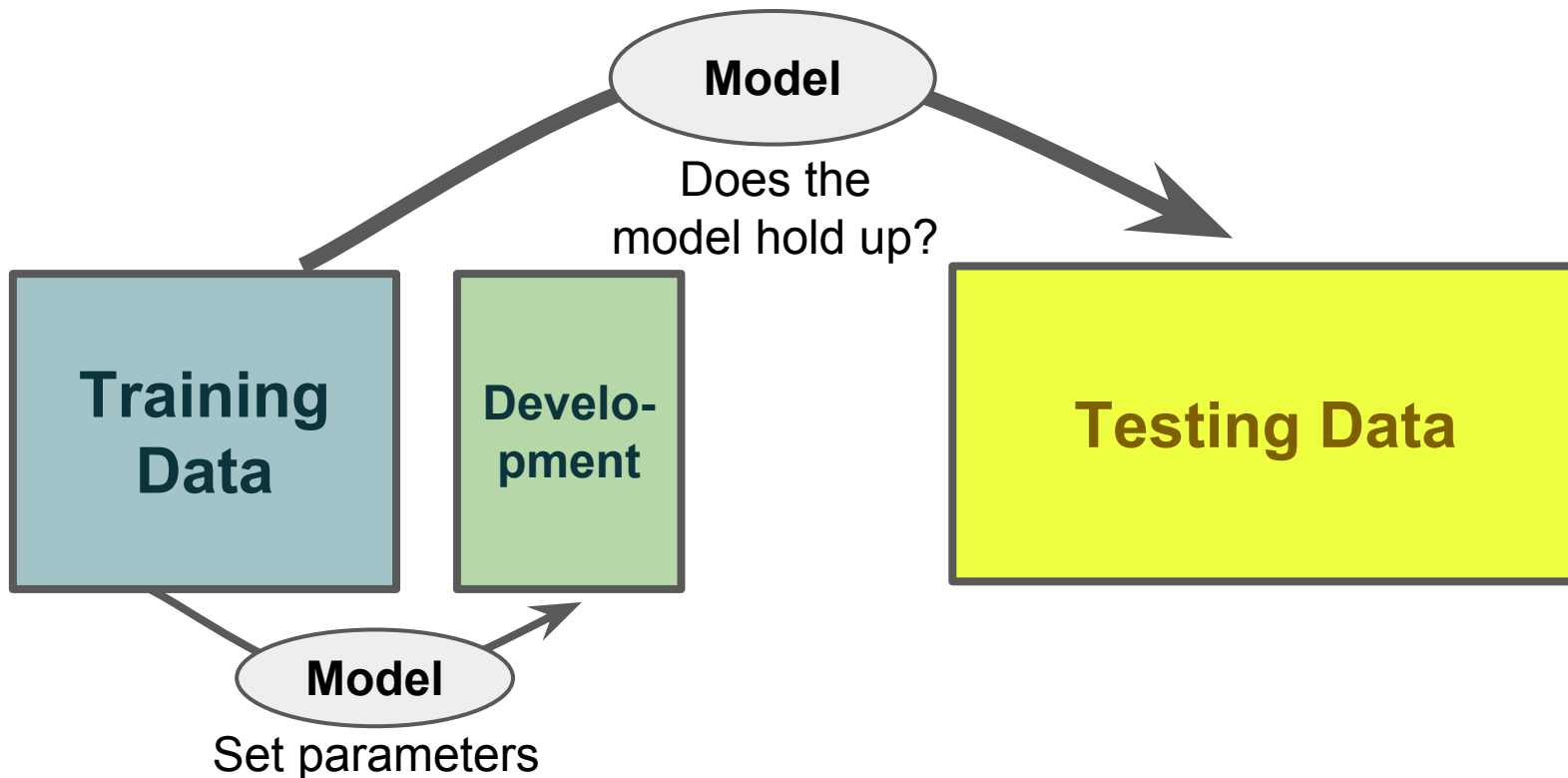


No closed form matrix solution, but often solved with coordinate descent.

$\lambda \|\beta\|_1$

Application:  $p \approx n$  or  $p \gg n$  (p: features; n: observations)

# Common Goal: Generalize to new data



# N-Fold Cross-Validation

Goal: Decent estimate of model accuracy



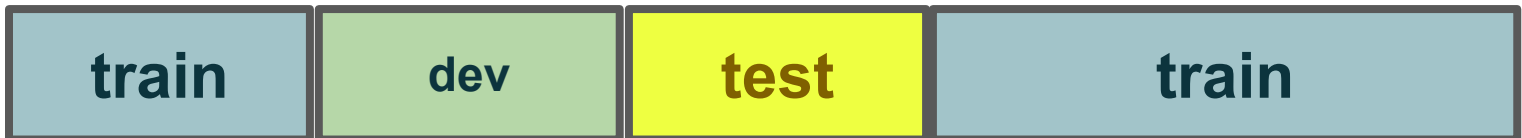
Iter 1



Iter 2



Iter 3



....

...

# Options for Distributing ML

1. **Distribute copies of entire dataset**
  - a. Train over all with different hyperparameters
  - b. Train different folds per worker node.

# Options for Distributing ML

1. **Distribute copies of entire dataset**
  - a. Train over all with different parameters
  - b. Train different folds per worker node.
  
2. **Distribute data**
  - a. Each node finds parameters for subset of data
  - b. Needs mechanism for updating parameters
    - i. Centralized parameter server
    - ii. Distributed All-Reduce

# Options for Distributing ML

## 1. **Distribute copies of entire dataset**

- a. Train over all with different parameters
- b. Train different folds per worker node.

## 2. **Distribute data**

- a. Each node finds parameters for subset of data
- b. Needs mechanism for updating parameters
  - i. Centralized parameter server
  - ii. Distributed All-Reduce

## 3. **Distribute model or individual operations** (e.g. matrix multiply)

# Options for Distributing ML

## 1. Distribute copies of entire dataset

- a. Train over all with different parameters
- b. Train different folds per worker node.

Pro: Easy; Good for compute-bound; Con: Requires data fit in worker memories

## 2. Distribute data

- a. Each node finds parameters for subset of data
- b. Needs mechanism for updating parameters
  - i. Centralized parameter server
  - ii. Distributed All-Reduce

Pro: Flexible to all situations; Con: Optimizing for subset is suboptimal

## 3. Distribute model or individual operations (e.g. matrix multiply)

Pro: Parameters can be localized Con: High communication for transferring Intermediar data.

# Options for Distributing ML

Done very often in practice. Not talked about much because it's mostly as easy as it sounds.

## 1. Distribute copies of entire dataset

- a. Train over all with different parameters
- b. Train different folds per worker node.

**Pro:** Easy; Good for compute-bound; **Con:** Requires data fit in worker memories

## 2. Distribute data

- a. Each node finds parameters for subset of data
- b. Needs mechanism for updating parameters
  - i. Centralized parameter server
  - ii. Distributed All-Reduce

**Pro:** Flexible to all situations; **Con:** Optimizing for subset is suboptimal

## 3. Distribute model or individual operations (e.g. matrix multiply)

**Pro:** Parameters can be localized **Con:** High communication for transferring Intermediar data.



# Options for Distributing ML

## 1. Distribute copies of entire dataset

- a. Train over all with different parameters
- b. Train different folds per worker node.

Done very often in practice. Not talked about much because it's mostly as easy as it sounds.

**Pro:** Easy; Good for compute-bound; **Con:** Requires data fit in worker memories

## 2. Distribute data

- a. Each node finds parameters for subset of data
- b. Needs mechanism for updating parameters
  - i. Centralized parameter server
  - ii. Distributed All-Reduce

Preferred method for big data or very complex models (i.e. models with many internal parameters).

**Pro:** Flexible to all situations; **Con:** Optimizing for subset is suboptimal

## 3. Distribute model or individual operations (e.g. matrix multiply)

**Pro:** Parameters can be localized **Con:** High communication for transferring Intermediar data.

# Options for Distributing ML

## 1. Distribute copies of entire dataset

- a. Train over all with different parameters
- b. Train different folds per worker node.

Done very often in practice. Not talked about much because it's mostly as easy as it sounds.

Pro: Easy; Good for compute-bound; Con: Requires data fit in worker memories

## 2. Distribute data

- a. Each node finds parameters for subset of data
- b. Needs mechanism for updating parameters
  - i. Centralized parameter server
  - ii. Distributed All-Reduce

### Data Parallelism

Preferred method for big data or very complex models (i.e. models with many internal parameters).

Pro: Flexible to all situations; Con: Optimizing for subset is suboptimal

## 3. Distribute model or individual operations (e.g. matrix multiply)

Pro: Parallelism can be leveraged

### Model Parallelism

Con: High communication for transferring Intermediar data.

# Cluster Distribution

Model Parallelism

Multiple devices on multiple machines

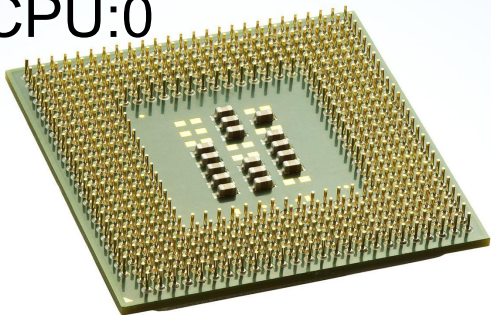
```
with tf.device("/cpu:1")  
    beta=tf.Variable(...)
```

```
with tf.device("/gpu:0")  
    y_pred=tf.matmul(beta,X)
```

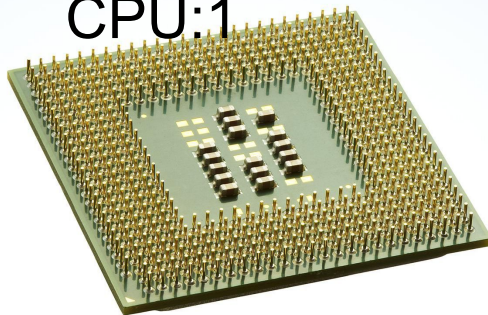
Transfer Tensors

Machine A

CPU:0

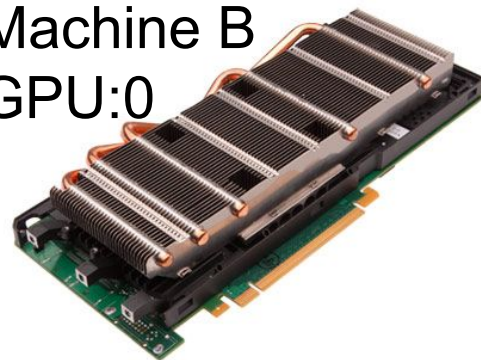


CPU:1



Machine B

GPU:0



# Cluster Distribution

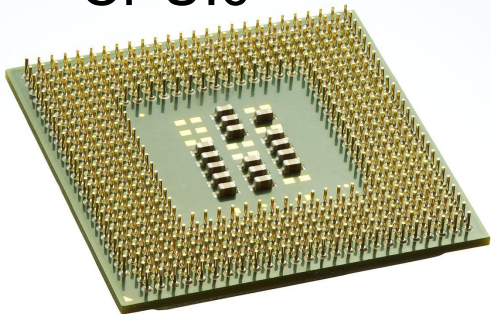
Data Parallelism

```
...  
beta=tf.Variable(...)  
pred=tf.matmul(beta,X)
```

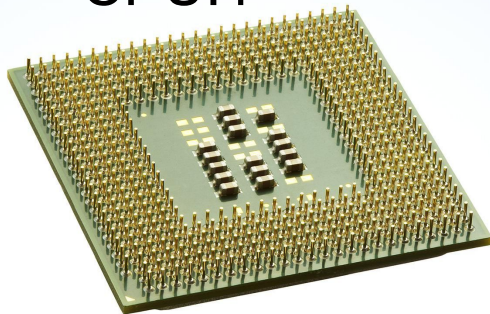
```
...  
beta=tf.Variable(...)  
pred=tf.matmul(beta,X)
```

```
...  
beta=tf.Variable(...)  
pred=tf.matmul(beta,X)
```

CPU:0



CPU:1



GPU:0



# Cluster Distribution

## Data Parallelism

```
...  
beta=tf.Variable(...)  
pred=tf.matmul(beta,X)
```

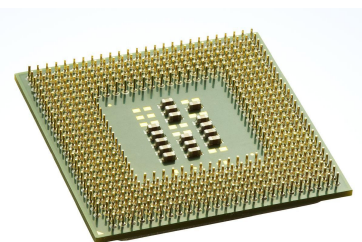
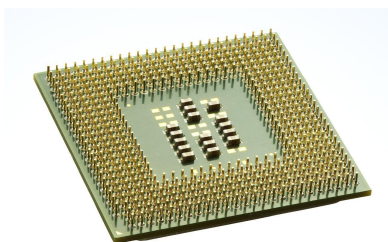
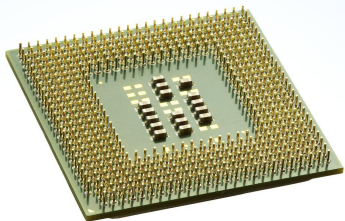
```
...  
beta=tf.Variable(...)  
pred=tf.matmul(beta,X)
```

```
...  
beta=tf.Variable(...)  
pred=tf.matmul(beta,X)
```

worker:0

worker:1

worker:2



# Distributing data

$X$

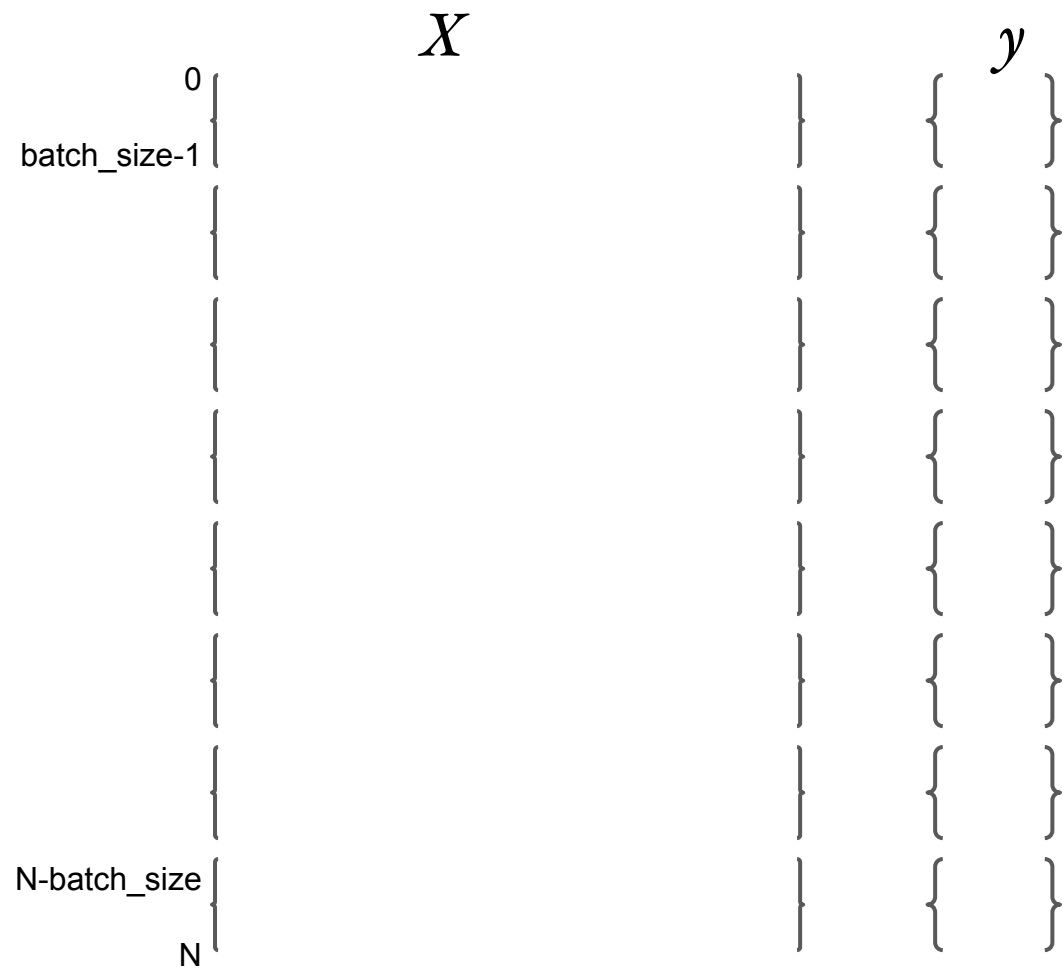
0  
N

}

$y$

}

# Distributing data

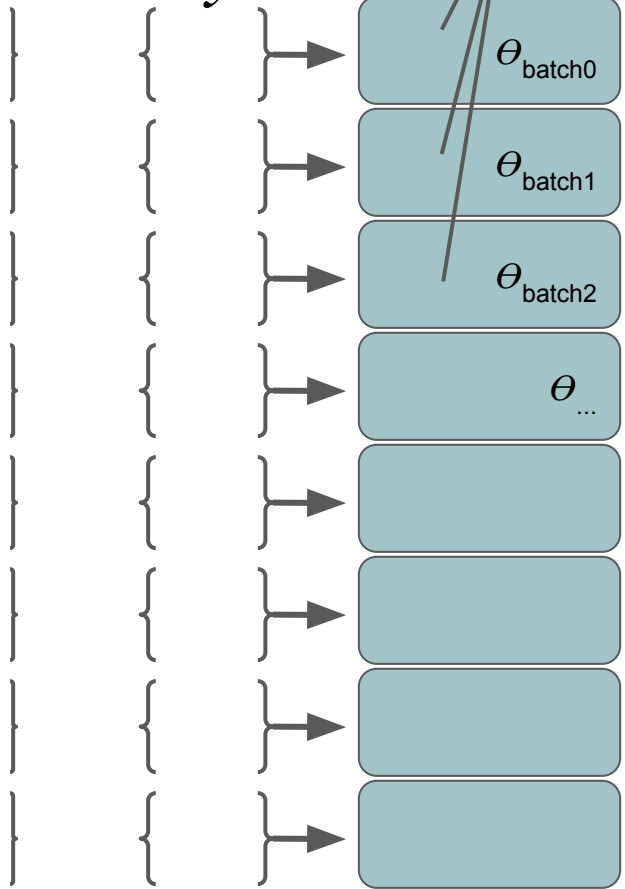
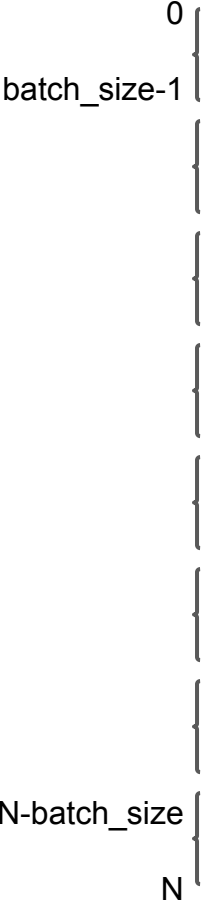


# Distributing data

learn parameters (i.e. weights),  
given graph with cost function  
and optimizer

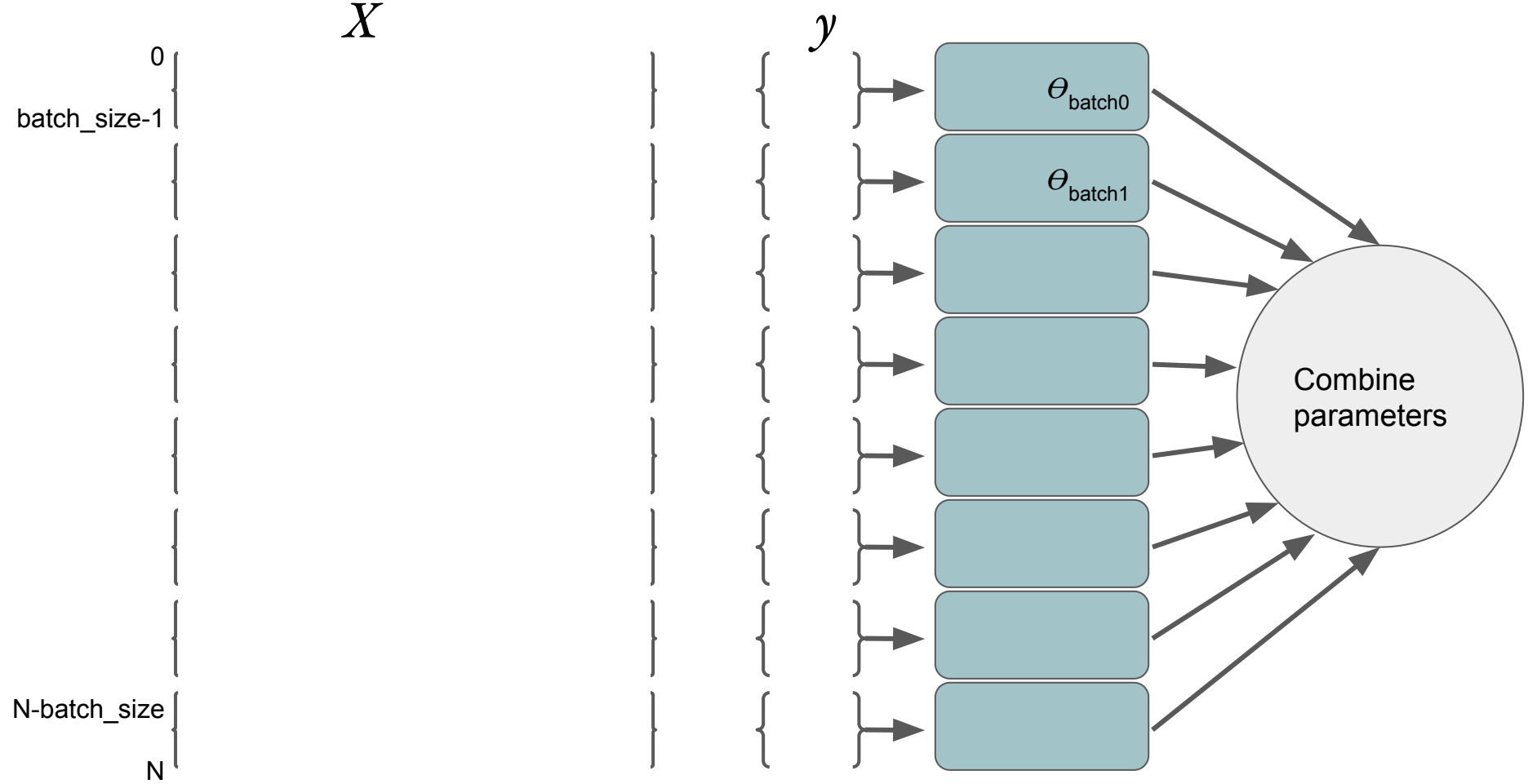
$X$

$y$





# Distributing data

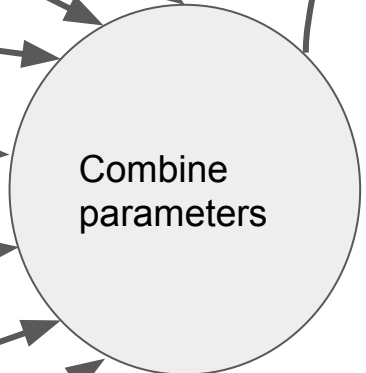
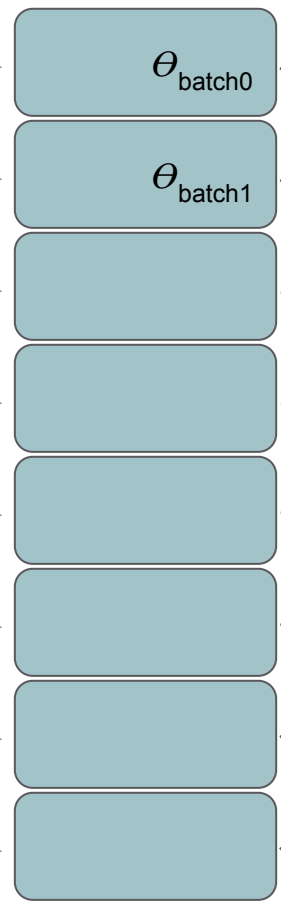


0  
batch\_size-1  
N-batch\_size  
N

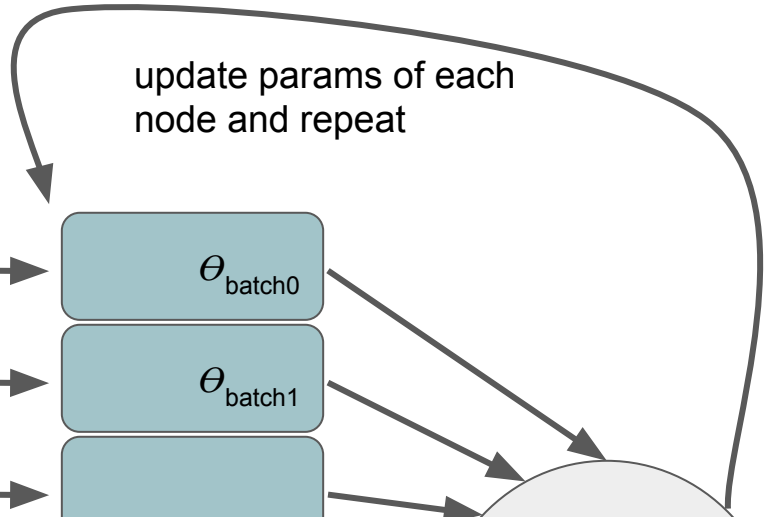
$X$



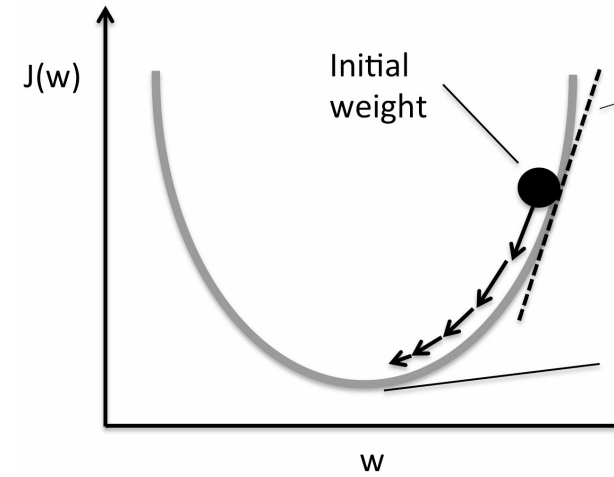
$y$



update params of each node and repeat



# Gradient Descent Options for Linear Regression



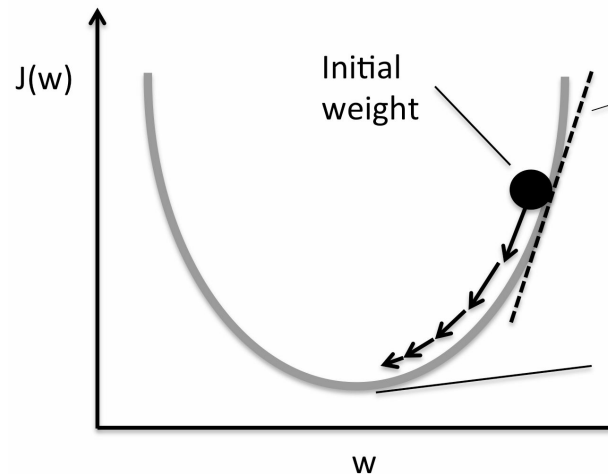
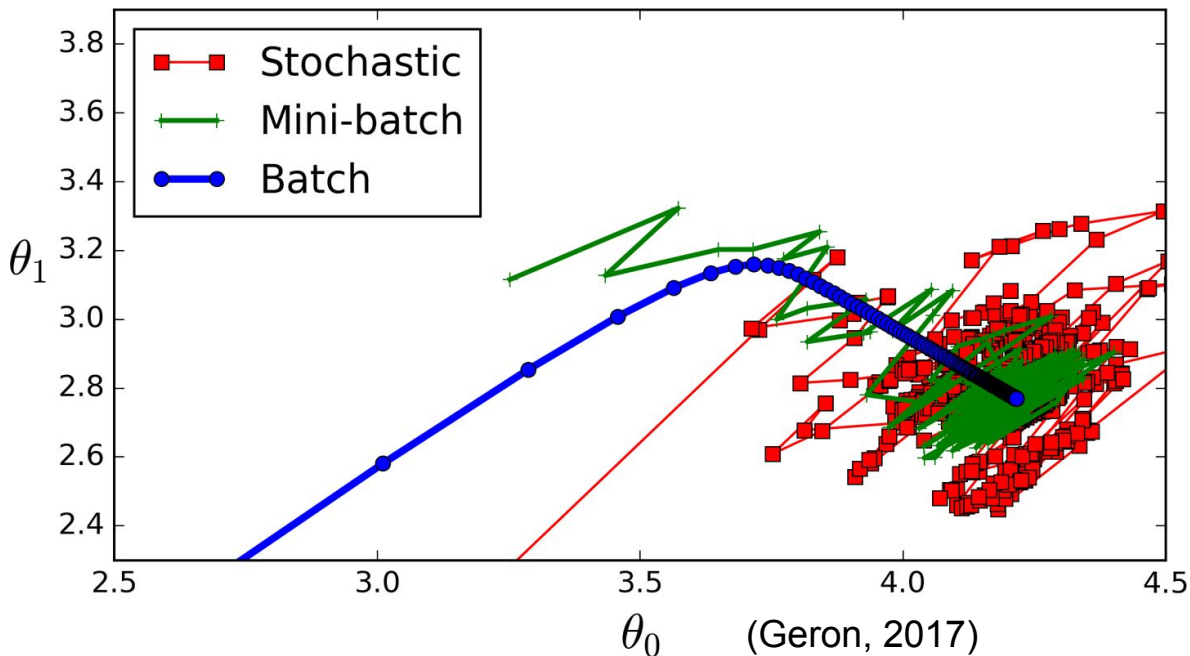
(Geron, 2017)

# Gradient Descent Options for Linear Regression

Batch Gradient Descent

Stochastic Gradient Descent: One example at a time

Mini-batch Gradient Descent:  $k$  examples at a time.

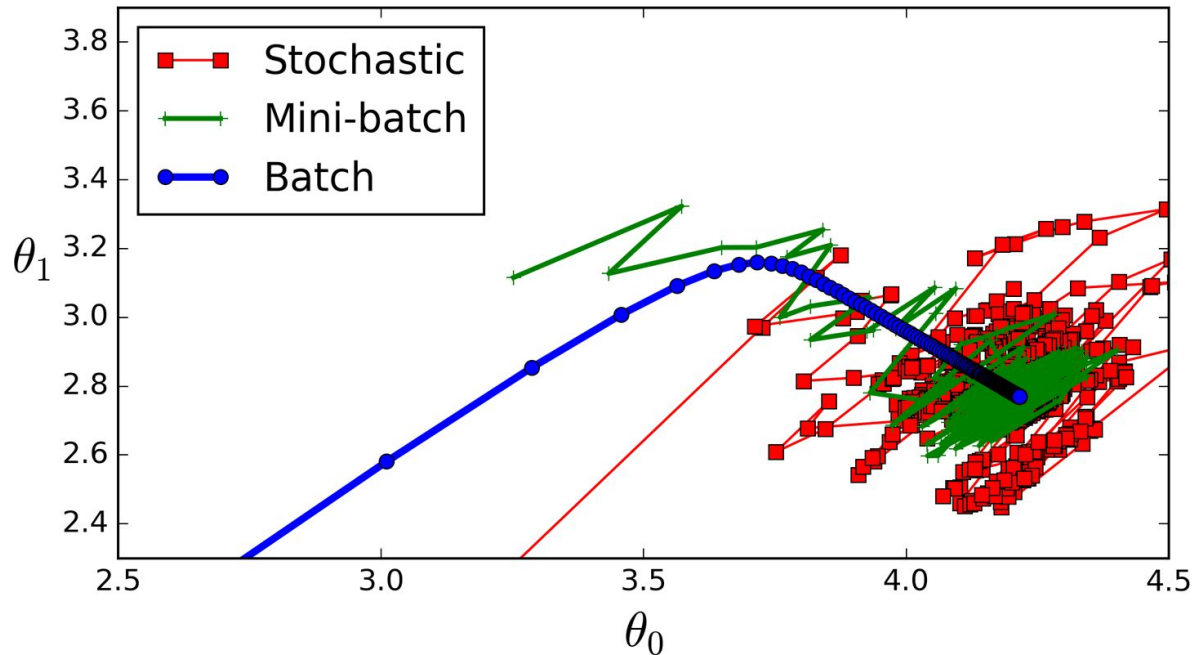


# Gradient Descent Options for Linear Regression

Batch Gradient Descent

Stochastic Gradient Descent: One example at a time

Mini-batch Gradient Descent: k examples at a time.



(Geron, 2017)

# From linear regression to neural networks

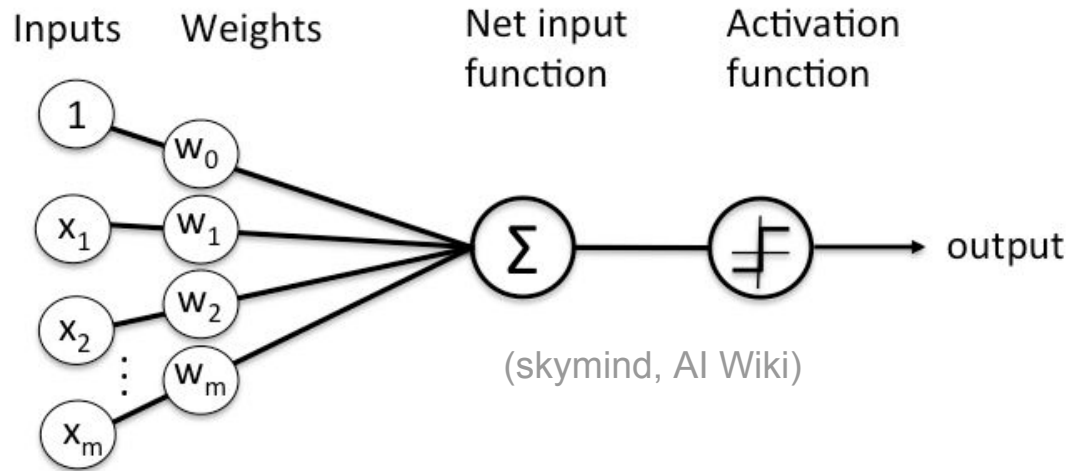
Linear Regression:  $y = wX$

Neural Network Nodes:  $output = f(wX)$

# From linear regression to neural networks

Linear Regression:  $y = wX$

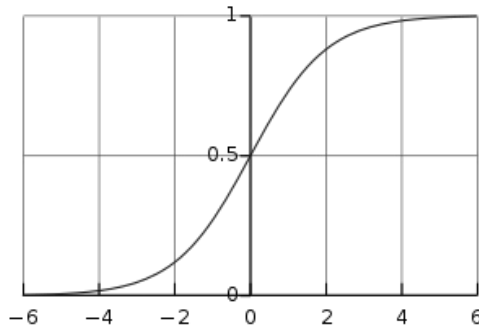
Neural Network Nodes:  $output = f(wX)$



# Common Activation Functions

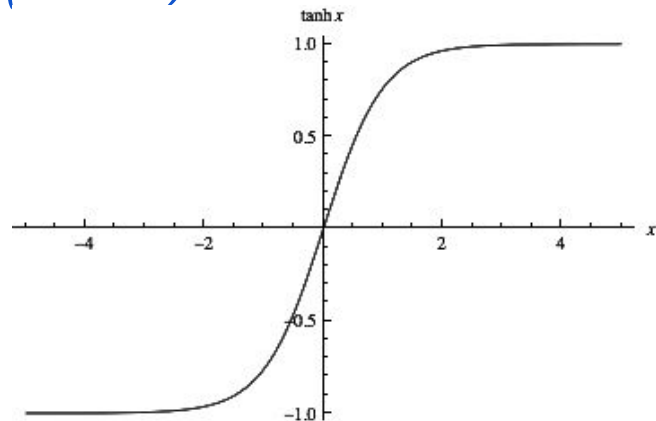
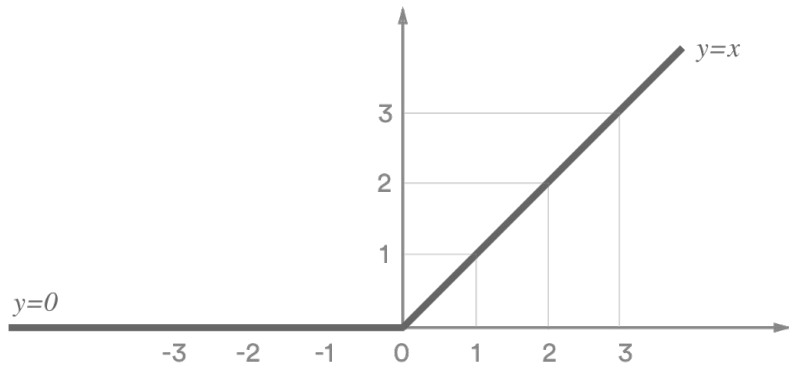
$$z = wX$$

Logistic:  $\sigma(z) = 1 / (1 + e^{-z})$



Hyperbolic tangent:  $\tanh(z) = 2\sigma(2z) - 1 = (e^{2z} - 1) / (e^{2z} + 1)$

Rectified linear unit (ReLU):  $ReLU(z) = \max(0, z)$

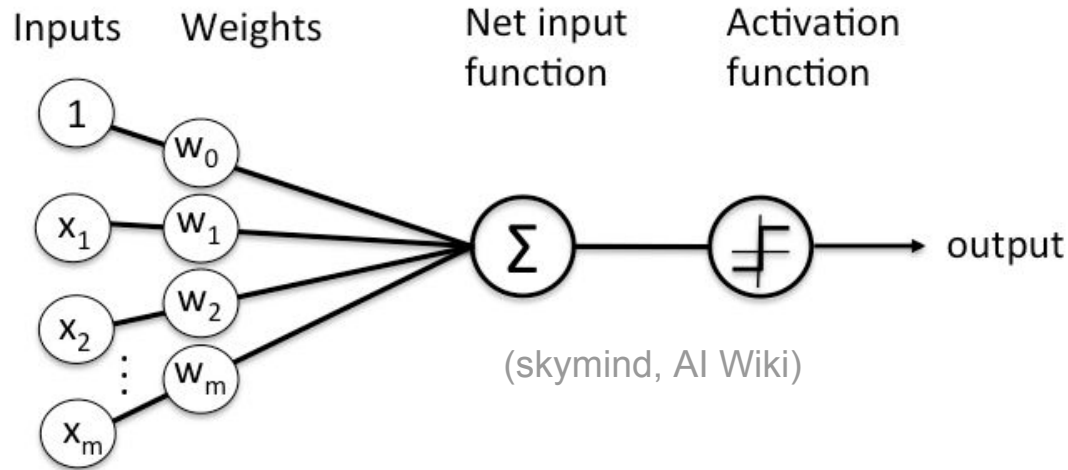




# From linear regression to neural networks

Linear Regression:  $y = wX$

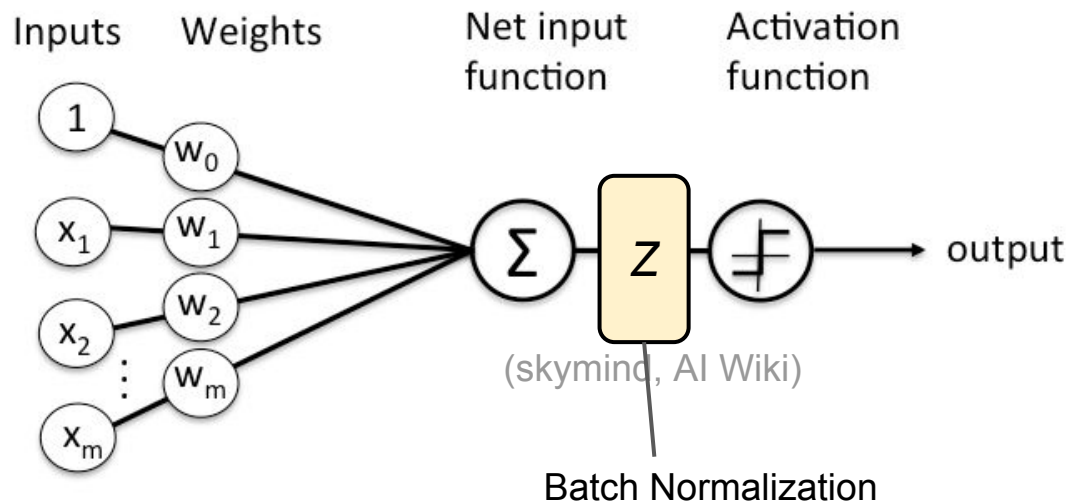
Neural Network Nodes:  $output = f(wX)$



# From linear regression to neural networks

Linear Regression:  $y = wX$

Neural Network Nodes:  $output = f(wX)$



# Batch Normalization

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

# Batch Normalization

This is just standardizing!  
(but within the current batch of observations)

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

# Batch Normalization

This is just standardizing!  
(but within the current batch of observations)

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;  
Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

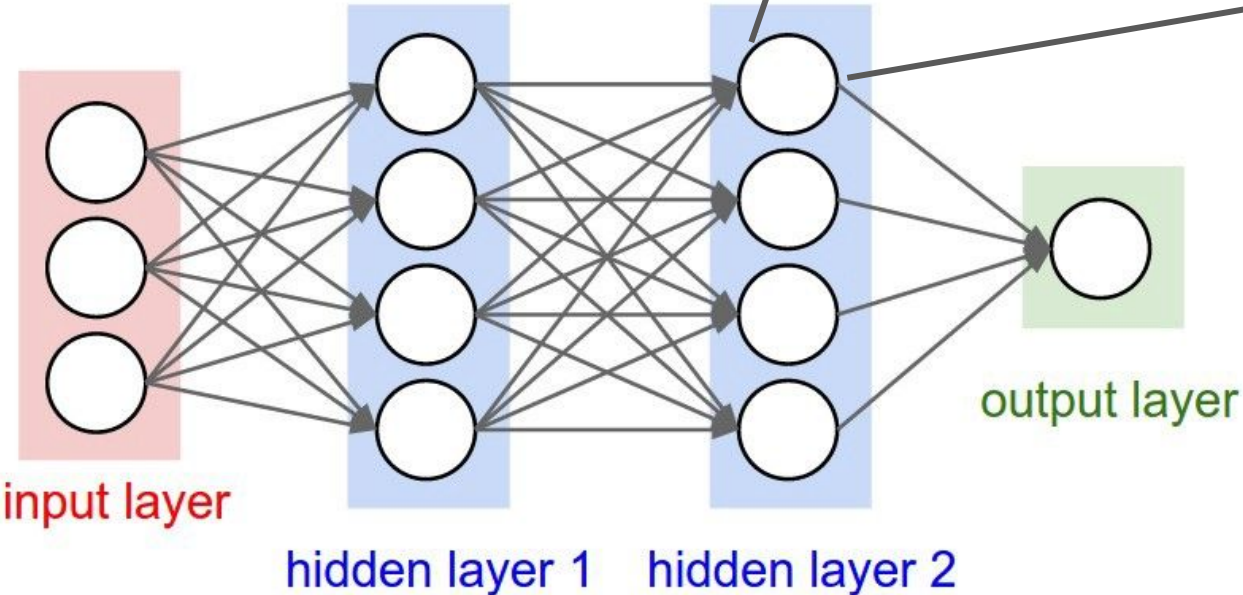
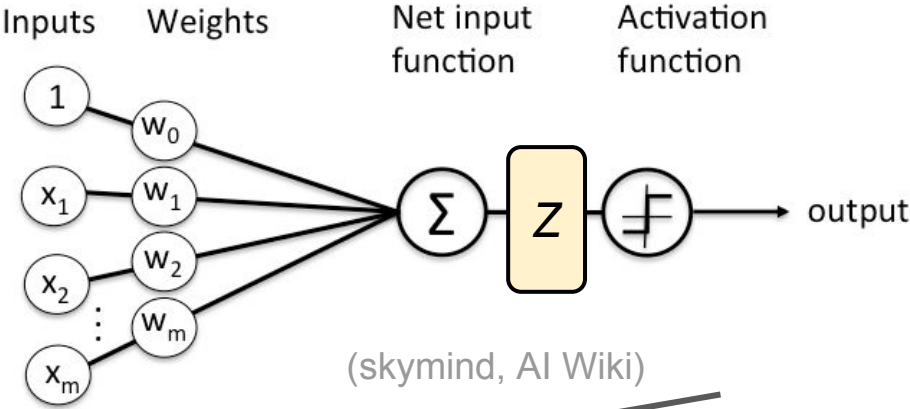
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

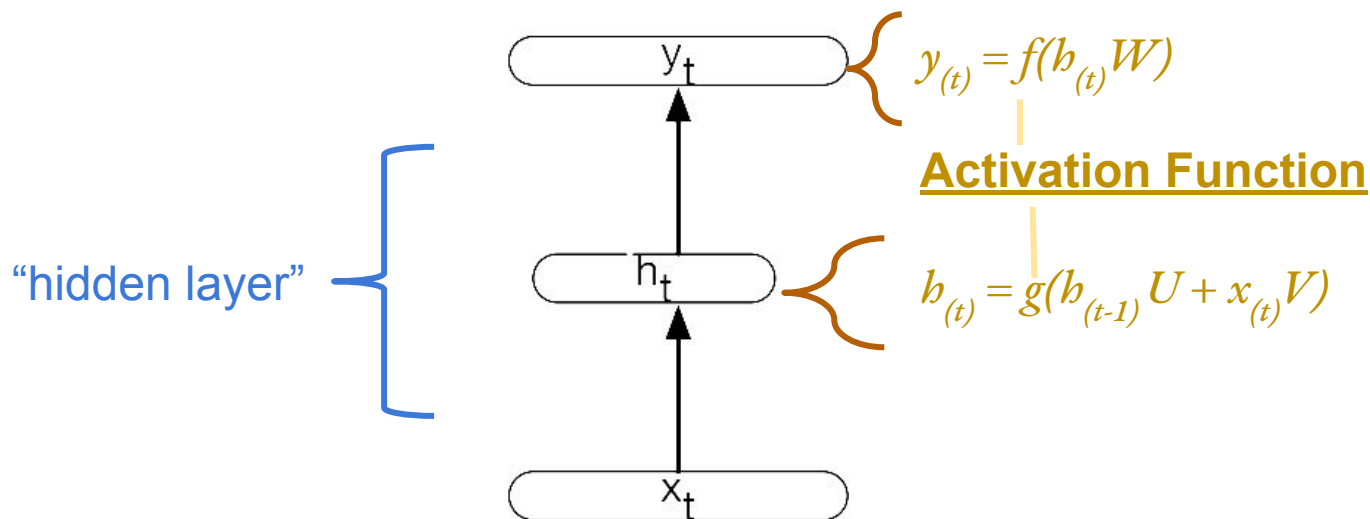
## Why?

- Empirically, it works!
- Conceptually, generally good for weight optimization to keep data within a reasonable range (dividing by sigma) and such that positive weights move it up and negative down (centering).
- Small effect: When done over mini-batches, adds regularization due to differences between batches.

# Feed Forward Network



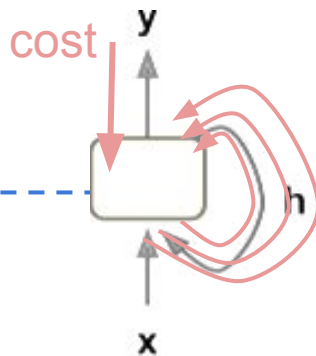
# Recurrent Neural Network



**Figure 9.2** Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep.

(Jurafsky, 2019)

# Optimization:



## Backward Propagation through Time

...

*#define forward pass graph:*

$h_{(0)} = \emptyset$

for  $i$  in range(1, len(x)):

$h_{(i)} = \text{tf.tanh}(\text{tf.matmul}(U, h_{(i-1)}) + \text{tf.matmul}(W, x_{(i)}))$  *#update hidden state*

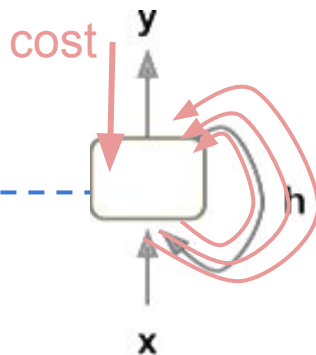
$y_{(i)} = \text{tf.softmax}(\text{tf.matmul}(V, h_{(i)}))$  *#update output*

...

$\text{cost} = \text{tf.reduce\_mean}(-\text{tf.reduce\_sum}(y * \text{tf.log}(y\_pred)))$



## Optimization:



## Backward Propagation through Time

...

```
#define forward pass graph:
```

```
 $h_{(0)} = \emptyset$ 
```

```
for i in range(1, len(x)):
```

```
     $h_{(i)} = \text{tf.tanh}(\text{tf.matmul}(U,$   
state
```

```
     $y_{(i)} = \text{tf.softmax}(\text{tf.matmul}$ 
```

```
...
```

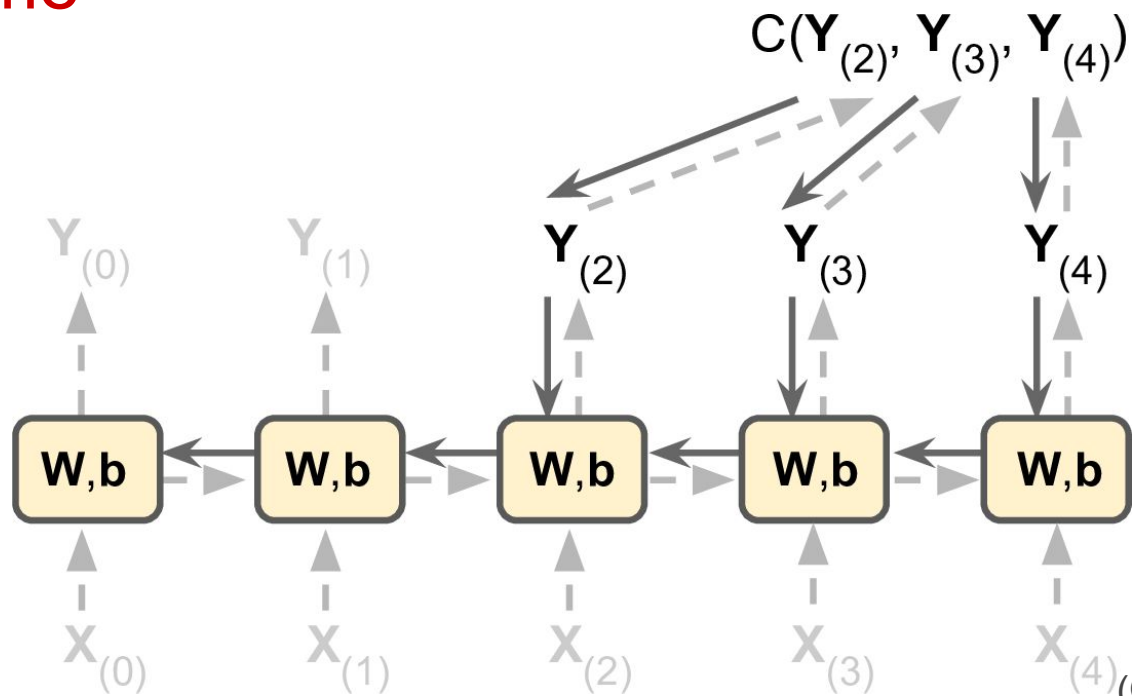
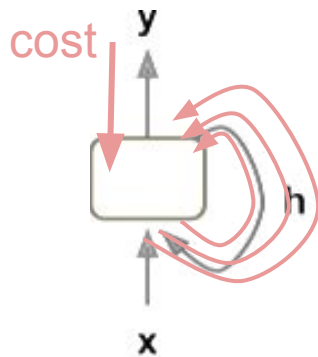
```
 $\text{cost} = \text{tf.reduce\_mean}(-\text{tf.reduce}$ 
```

To find the gradient for the overall graph, we use **back propogation**, which *essentially* chains together the gradients for each node (function) in the graph.

With many recursions, the gradients can vanish or explode (become too large or small for floating point operations).

Optimization:

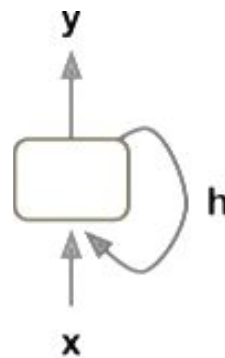
Backward Propagation  
through Time



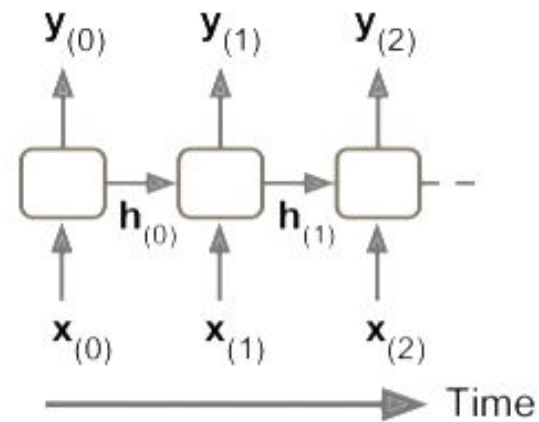
(Geron, 2017)

# How to address exploding and vanishing gradients?

Dominant approach: Use Long Short Term Memory Networks (LSTM)



RNN model

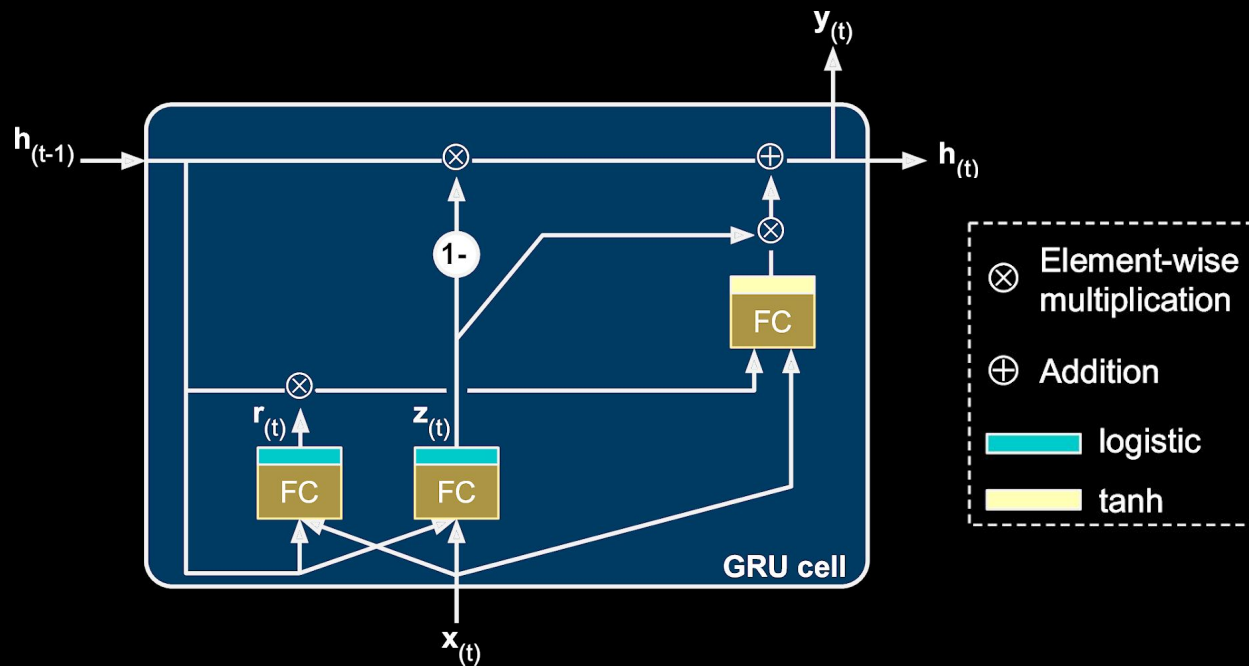


“unrolled” depiction

(Geron, 2017)

# The GRU Cell

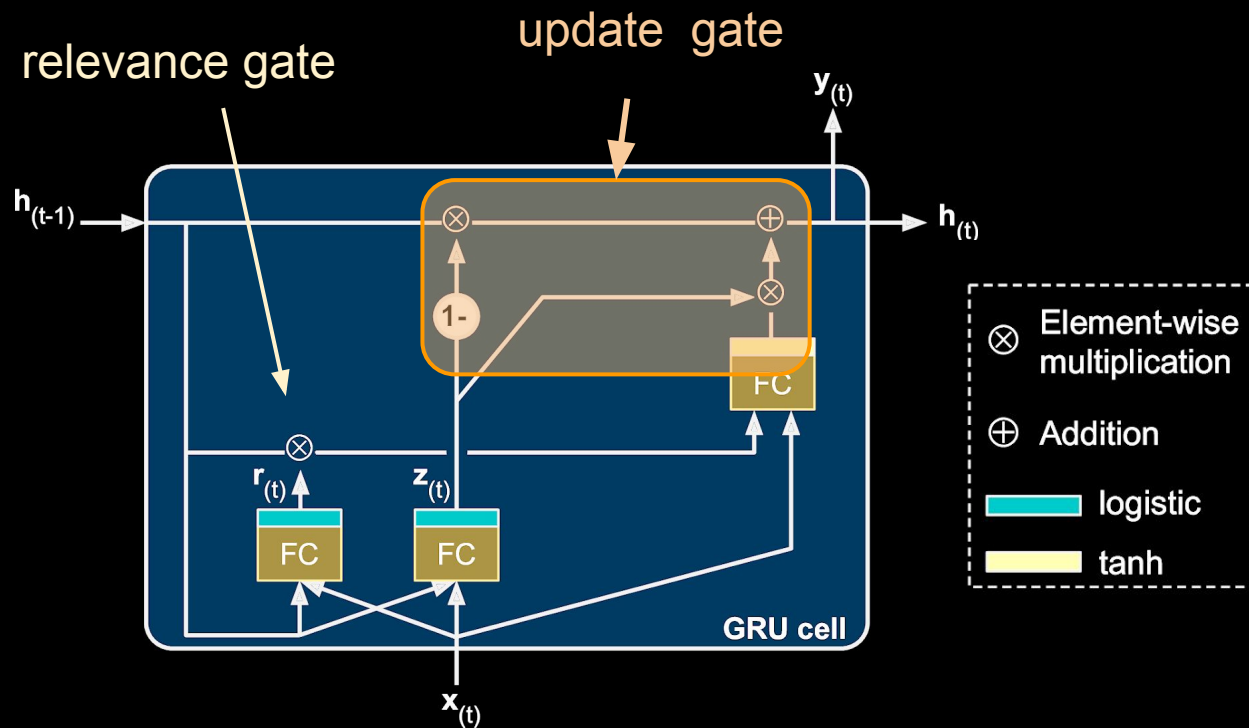
## Gated Recurrent Unit



(Geron, 2017)

# The GRU

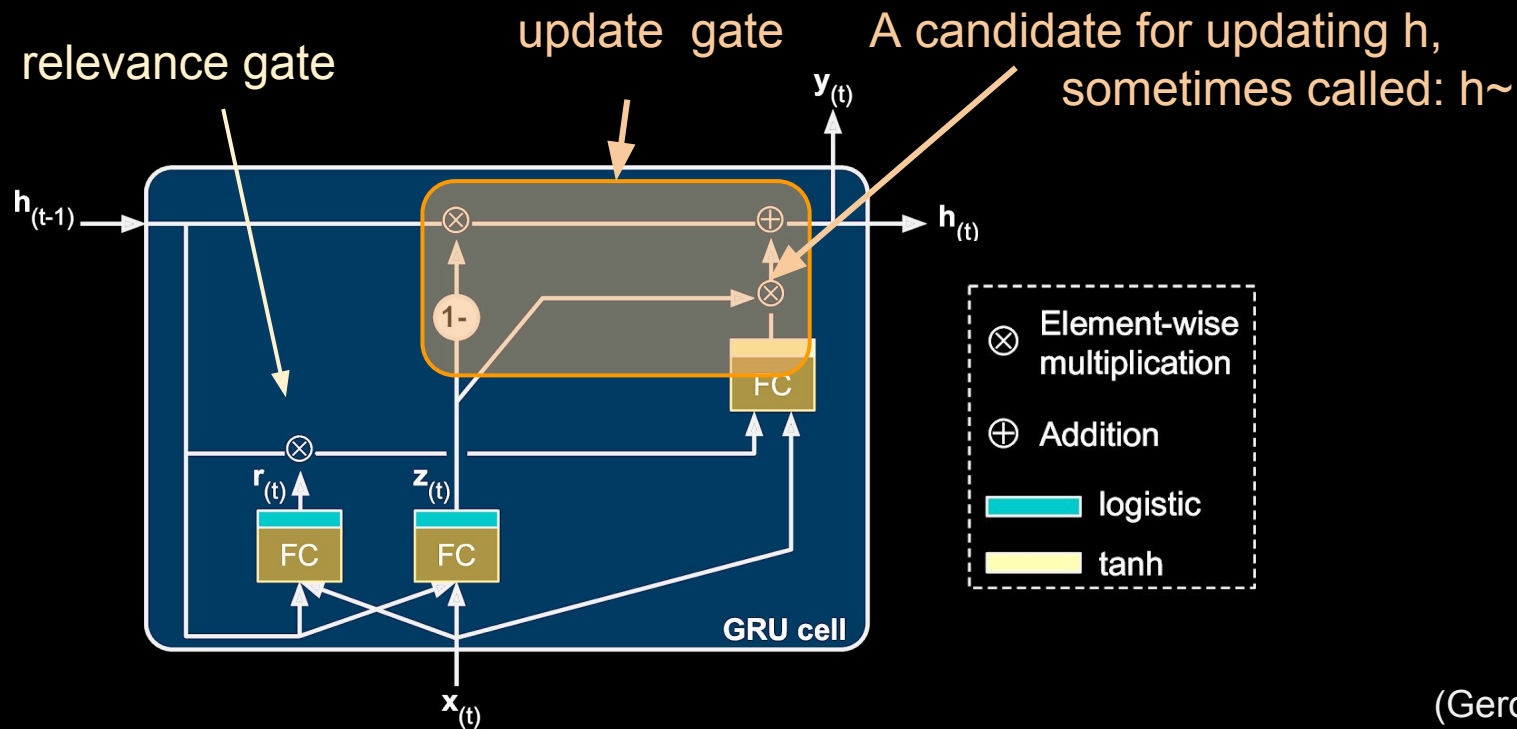
## Gated Recurrent Unit



(Geron, 2017)

# The GRU

## Gated Recurrent Unit



(Geron, 2017)

# The GRU

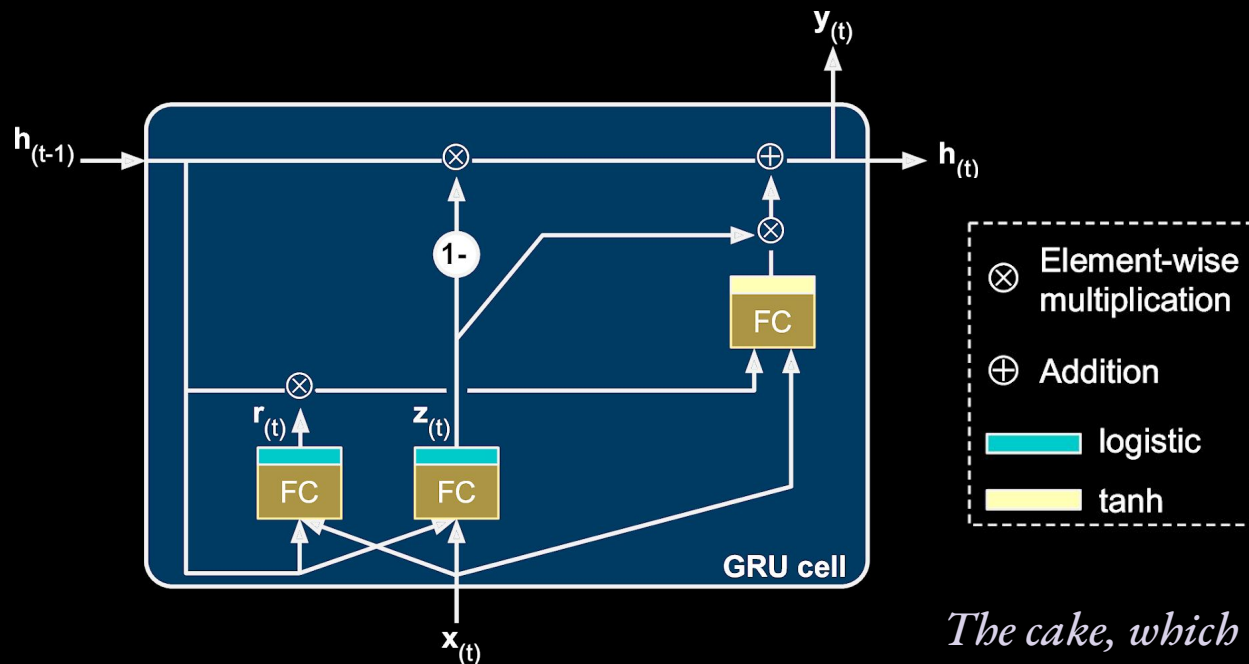
## Gated Recurrent Unit

$$\mathbf{z}_{(t)} = \sigma(\mathbf{W}_{xz}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hz}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_z)$$

$$\mathbf{r}_{(t)} = \sigma(\mathbf{W}_{xr}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hr}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_r)$$

$$\mathbf{g}_{(t)} = \tanh(\mathbf{W}_{xg}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^T \cdot (\mathbf{r}_{(t)} \otimes \mathbf{h}_{(t-1)}) + \mathbf{b}_g)$$

$$\mathbf{h}_{(t)} = \mathbf{z}_{(t)} \otimes \mathbf{h}_{(t-1)} + (1 - \mathbf{z}_{(t)}) \otimes \mathbf{g}_{(t)}$$



*The cake, which contained candles, was eaten.*

# What about the gradient?

$$\mathbf{z}_{(t)} = \sigma(\mathbf{W}_{xz}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hz}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_z)$$

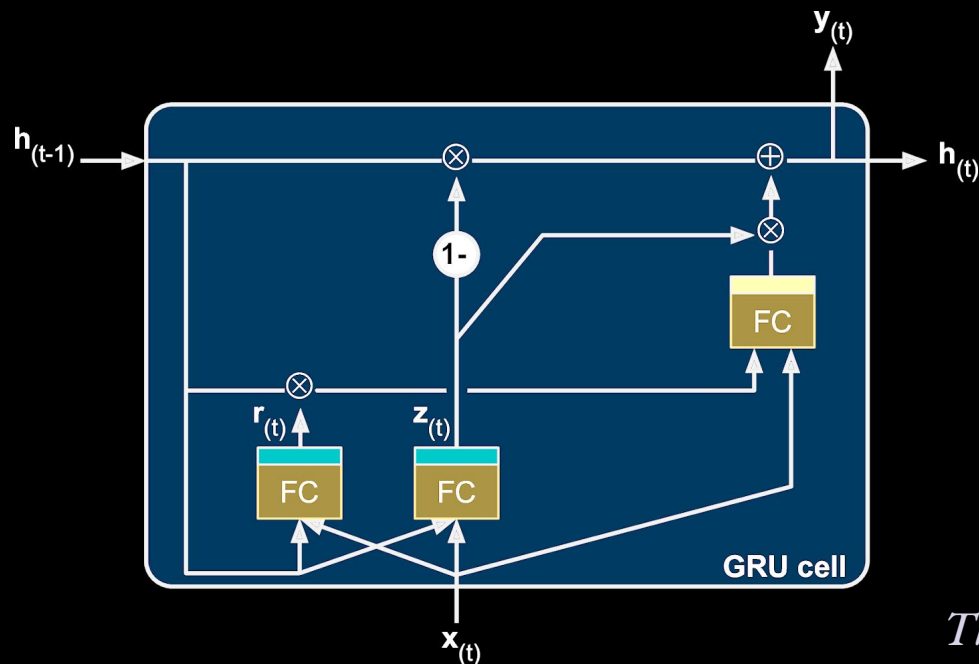
$$\mathbf{r}_{(t)} = \sigma(\mathbf{W}_{xr}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hr}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_r)$$

$$\mathbf{g}_{(t)} = \tanh(\mathbf{W}_{xg}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^T \cdot (\mathbf{r}_{(t)} \otimes \mathbf{h}_{(t-1)}) + \mathbf{b}_g)$$

$$\mathbf{h}_{(t)} = \mathbf{z}_{(t)} \otimes \mathbf{h}_{(t-1)} + (1 - \mathbf{z}_{(t)}) \otimes \mathbf{g}_{(t)}$$

The gates (i.e. multiplications based on a logistic) often end up keeping the hidden state exactly (or nearly exactly) as it was. Thus, for most dimensions of  $\mathbf{h}$ ,

$$\mathbf{h}_{(t)} \approx \mathbf{h}_{(t-1)}$$



*The cake, which contained candles, was eaten.*



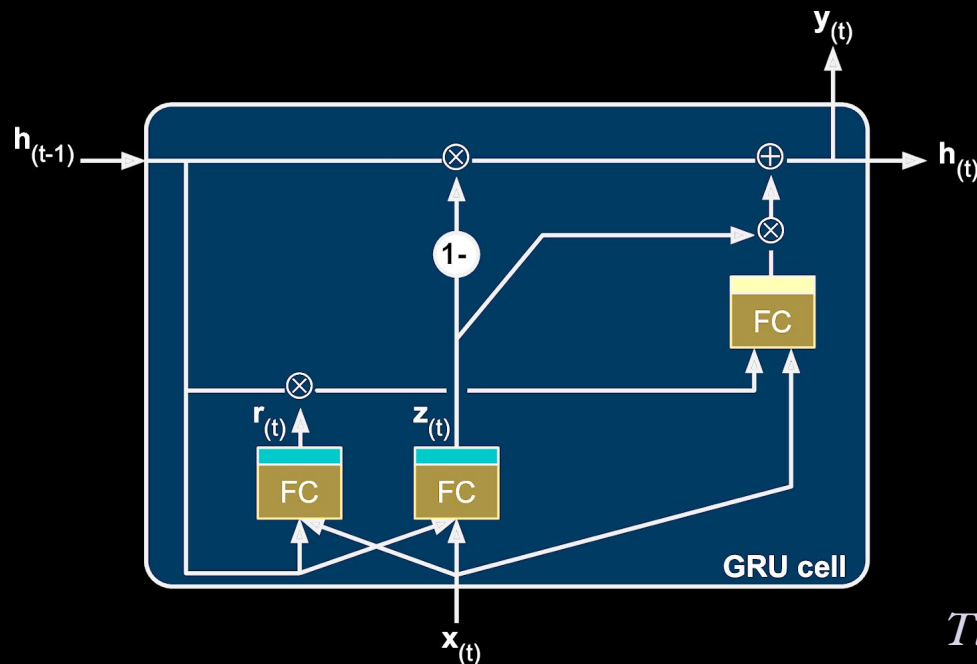
# What about the gradient?

$$\mathbf{z}_{(t)} = \sigma(\mathbf{W}_{xz}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hz}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_z)$$

$$\mathbf{r}_{(t)} = \sigma(\mathbf{W}_{xr}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hr}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_r)$$

$$\mathbf{g}_{(t)} = \tanh(\mathbf{W}_{xg}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^T \cdot (\mathbf{r}_{(t)} \otimes \mathbf{h}_{(t-1)}) + \mathbf{b}_g)$$

$$\mathbf{h}_{(t)} = \mathbf{z}_{(t)} \otimes \mathbf{h}_{(t-1)} + (1 - \mathbf{z}_{(t)}) \otimes \mathbf{g}_{(t)}$$



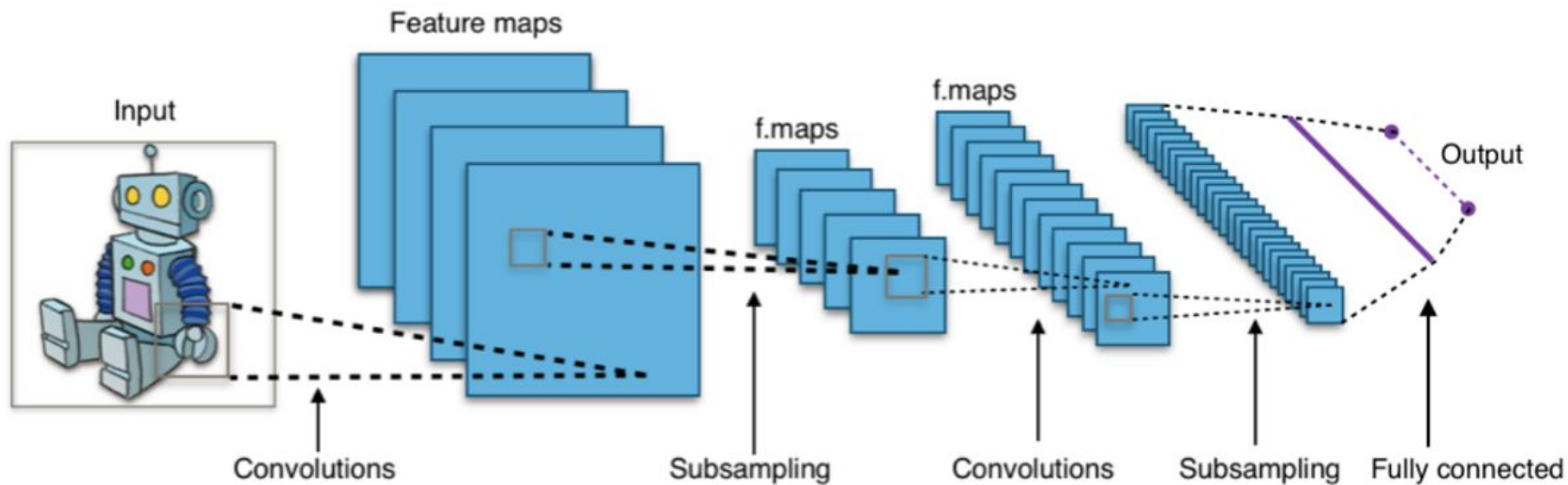
The gates (i.e. multiplications based on a logistic) often end up keeping the hidden state exactly (or nearly exactly) as it was. Thus, for most dimensions of  $\mathbf{h}$ ,

$$\mathbf{h}_{(t)} \approx \mathbf{h}_{(t-1)}$$

This tends to keep the gradient from vanishing since the same values will be present through multiple times in backpropagation through time. (The same idea applies to LSTMs but is easier to see here).

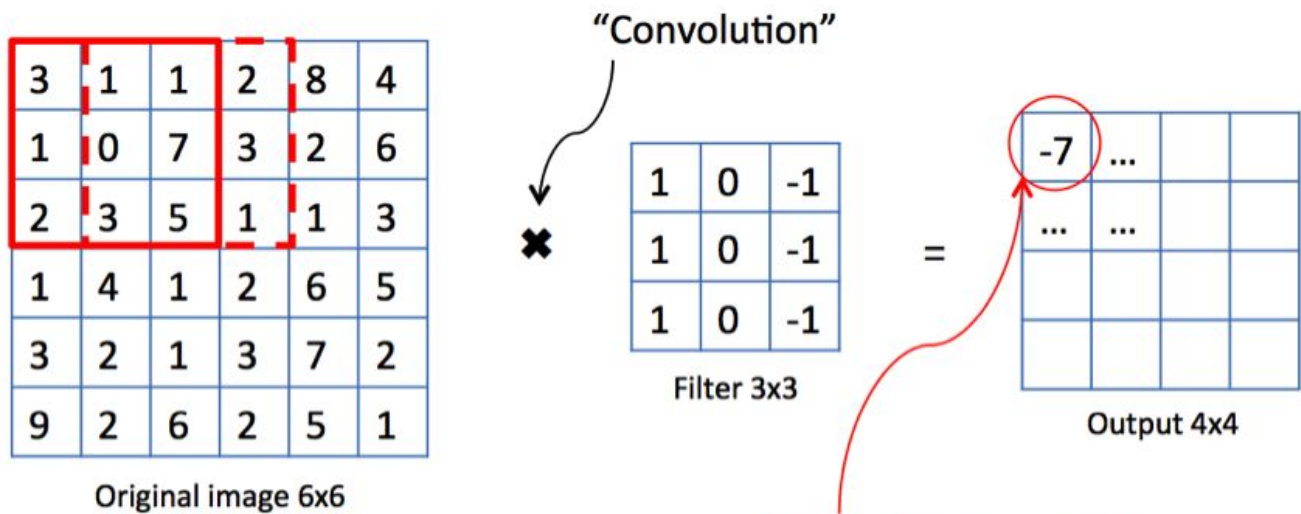
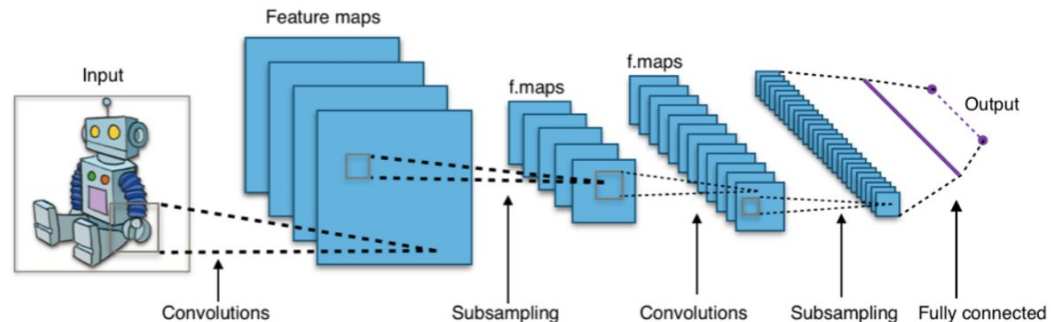
*The cake, which contained candles, was eaten.*

# Convolutional Layer



(wikipedia)

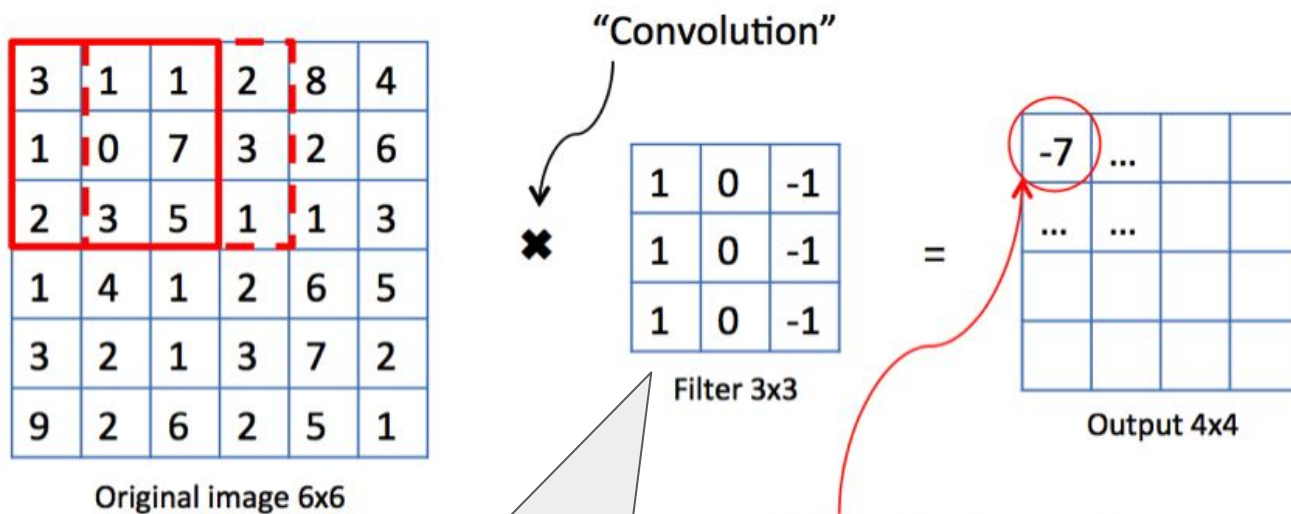
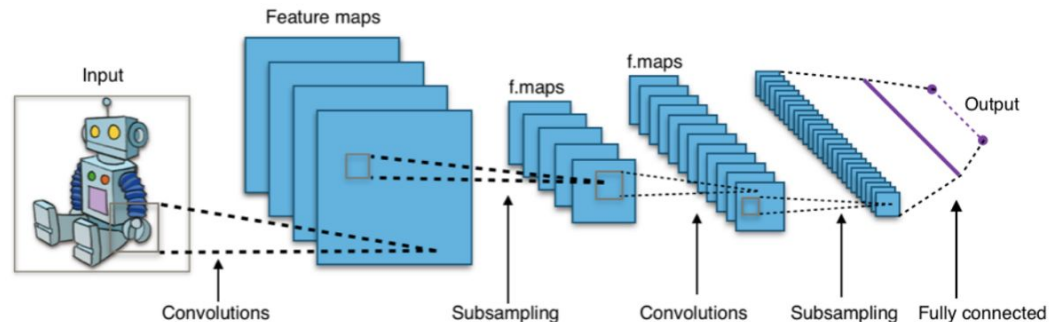
# Convolutional Layer



Result of the element-wise product and sum of the filter matrix and the original image

(Barter, 2018)

# Convolutional Layer



Breakthrough in image classification: Let the model automatically learn the filter weights!

Result of the element-wise product and sum of the filter matrix and the original image

# How to train deep models for classification?

Short Answer: Same as logistic regression.

```
RNN_cost = tf.reduce_mean(-tf.reduce_sum(y*tf.log(y_pred))  
#where did this come from?
```

Logistic Regression Likelihood:  $L(\beta_0, \beta_1, \dots, \beta_k | X, Y) = \prod_{i=1}^N p(x_i)^{y_i} (1 - p(x_i))^{1-y_i}$

Log Likelihood:  $\ell(\beta) = \sum_{i=1}^N y_i \log p(x_i) + (1 - y_i) \log (1 - p(x_i))$

Log Loss:  $J(\beta) = -\frac{1}{N} \sum_{i=1}^N y_i \log p(x_i) + (1 - y_i) \log (1 - p(x_i))$

Cross-Entropy Cost:  $J = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{|V|} y_{i,j} \log p(x_{i,j})$  (a "multiclass" log loss)

Final Cost Function:  $J^{(t)} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{|V|} y_{i,j}^{(t)} \log \hat{y}_{i,j}^{(t)}$  -- "cross entropy error"

# Summary

- Goal is accurate prediction of  $y$  (outcome) given features ( $x$ )
- Use L1 or L2 penalization (as a regularization) to avoid overfit
- Reason for Train, Dev, Test split
- Components of a neural network
- Batch Normalization
- Distribution options: why is data parallelism preferred?
- Recurrent Neural Network
- Convolution Operation with Filters