# Pytorch Tutorial 3

# The perceptron
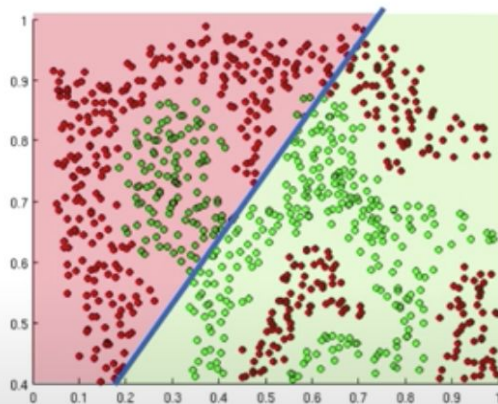
$$\hat{y} = g\left( w_0 + \boldsymbol{X}^T \boldsymbol{W} \right)$$



| Inputs | Weights | Sum | Non-Linearity | Output |

From: http://introtodeeplearning.com/slides/6S191_MIT_DeepLearning_L1.pdf

# Importance of Activation Functions

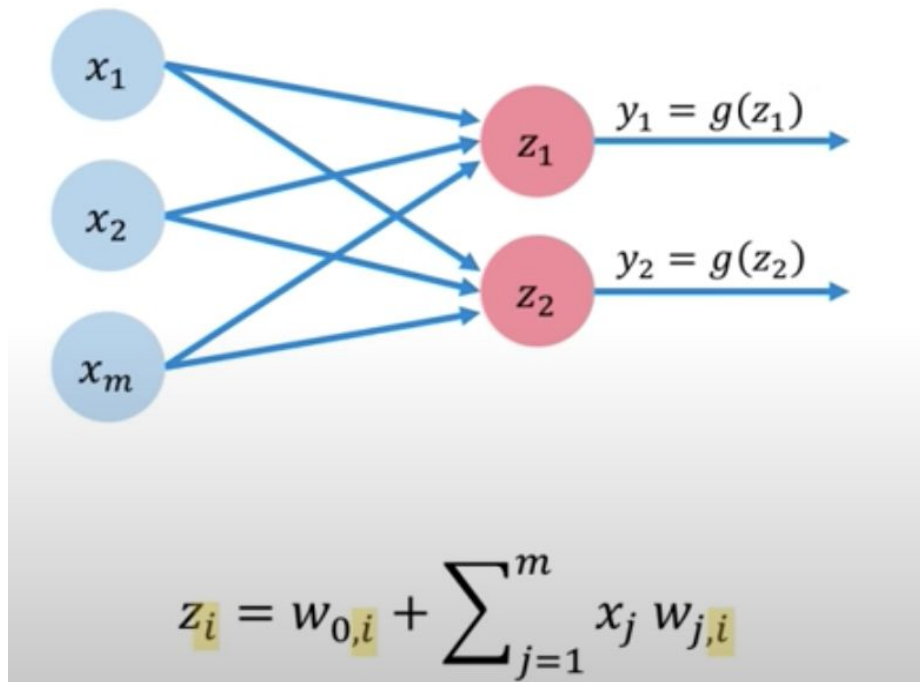*The purpose of activation functions is to **introduce non-linearities** into the network*



Linear activation functions produce linear decisions no matter the network size

Non-linearities allow us to approximate arbitrarily complex functions

From: http://introtodeeplearning.com/slides/6S191_MIT_DeepLearning_L1.pdf

# Multi Output Perceptron



$$z_i = w_{0,i} + \sum_{j=1}^{m} x_j \, w_{j,i}$$

From: http://introtodeeplearning.com/slides/6S191_MIT_DeepLearning_L1.pdf

# Single Layer Neural Network



$$z_i = w_{0,i}^{(1)} + \sum_{j=1}^{m} x_j \, w_{j,i}^{(1)} \qquad \hat{y}_i = g\left(w_{0,i}^{(2)} + \sum_{j=1}^{d_1} z_j \, w_{j,i}^{(2)}\right)$$

From: http://introtodeeplearning.com/slides/6S191_MIT_DeepLearning_L1.pdf

# Diving Deep into Supervised Training

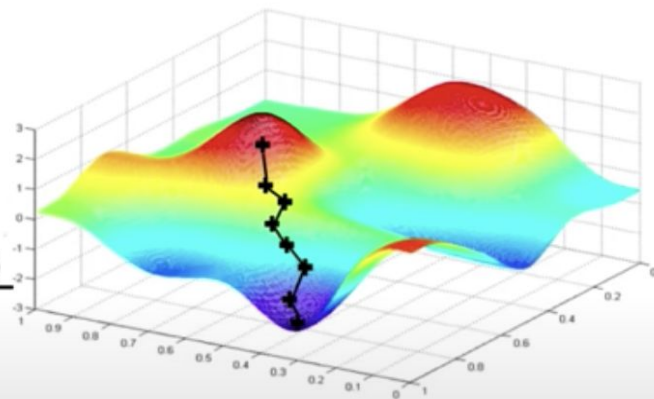★   Example: Supervised training for a perceptron and a binary classification

# Example: Supervised training for a perceptron and a binary classification

1) Choosing a model
   a) E.g. Perceptron
2) Choosing a loss function
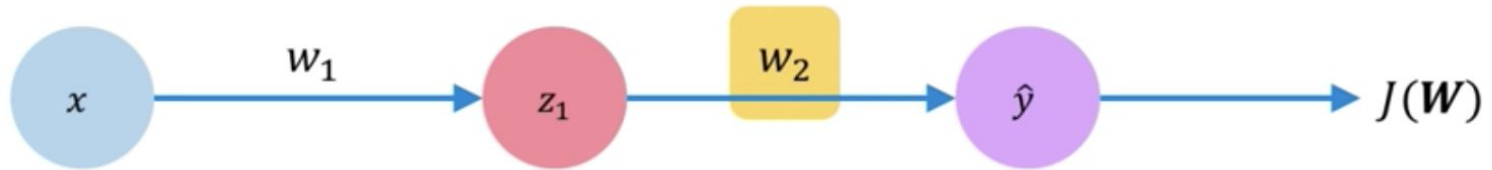3) Choosing an optimizer
   a) E.g. SGD, Adam

# Stochastic Gradient Descent

## Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$

2. Loop until convergence:

3.       Pick batch of $B$ data points

4.       Compute gradient, $\dfrac{\partial J(W)}{\partial W} = \dfrac{1}{B}\sum_{k=1}^{B} \dfrac{\partial J_k(W)}{\partial W}$

5.       Update weights, $W \leftarrow W - \eta\dfrac{\partial J(W)}{\partial W}$

6. Return weights



From: http://introtodeeplearning.com/slides/6S191_MIT_DeepLearning_L1.pdf
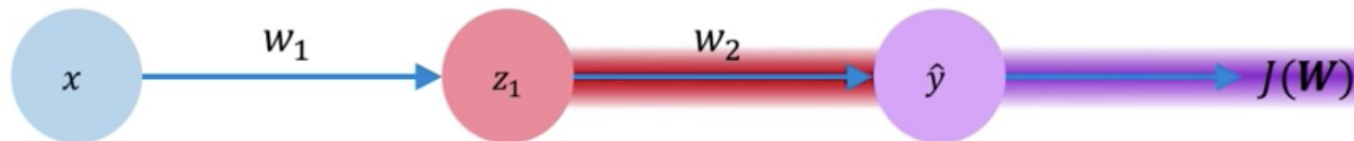
# Computing Gradients: Backpropagation



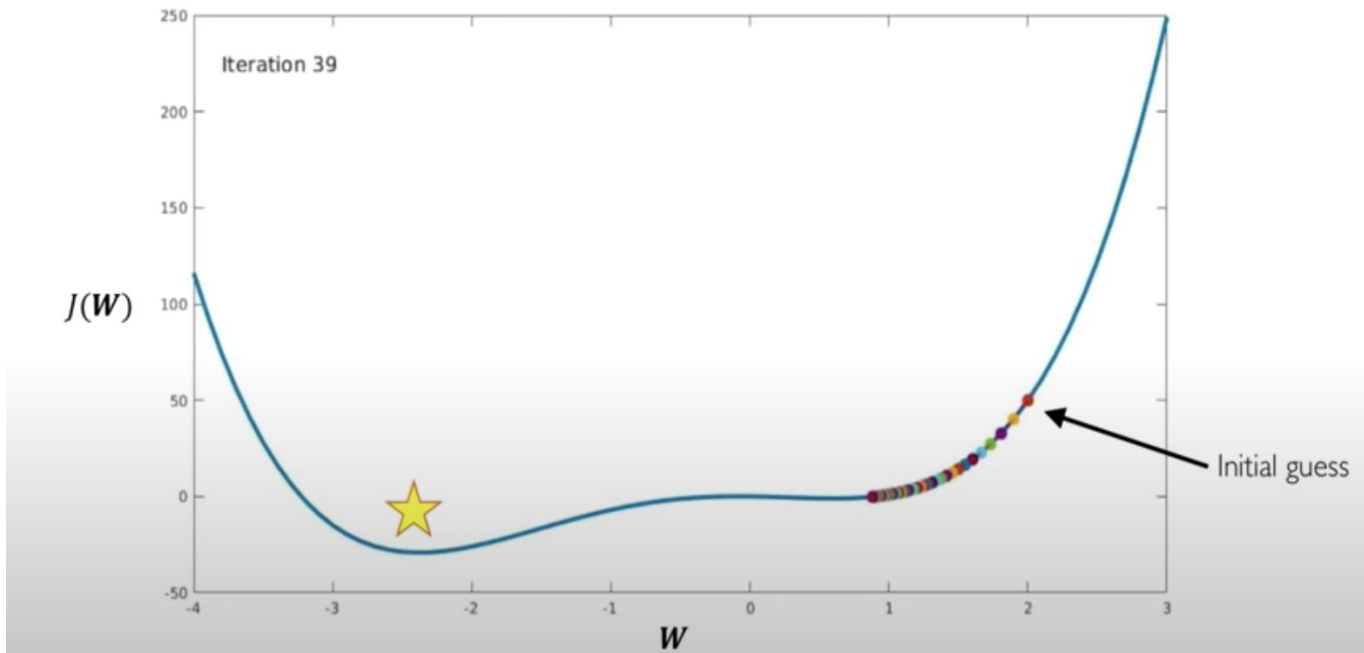*How does a small change in one weight (ex. $w_2$) affect the final loss $J(W)$?*

From:

# Computing Gradients: Backpropagation



$$\frac{\partial J(\boldsymbol{W})}{\partial w_2} = \frac{\partial J(\boldsymbol{W})}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w_2}$$

# Setting the Learning Rate

**Small learning rate** *converges slowly and gets stuck in false local minima*

# Setting the Learning Rate

*Large learning rates* overshoot, become unstable and diverge

# Adaptive Learning Rates

- Learning rates are no longer fixed
- Can be made larger or smaller depending on:
  - how large gradient is
  - how fast learning is happening
  - size of particular weights
  - etc...

E.g. Adam

# Example 3-10 Instantiating the Adam Optimizer

```python
import torch.nn as nn
import torch.optim as optim

input_dim = 2
lr = 0.001

perceptron = Perceptron(input_dim=input_dim)
bce_loss = nn.BCELoss()
optimizer = optim.Adam(params=perceptron.parameters(), lr=lr)
```

# Example 3-11 A supervised training loop for a perceptron and binary classification

```python
#each epoch is a complete pass over the training data
for epoch_i in range(n_epochs):
  #the inner loop is over the batches in the dataset
  for batch_i in range(n_batches):

    #Step 0: Get the data
    x_data, y_target = get_toy_data(batch_size)

    #Step 1: Clear the gradients
    perceptron.zero_grad()

    #Step 2:Compute the forward pass of the model
    y_pred = perceptron(x_data, apply_sigmoid=True)

    #Step 3: Compute the loss value that we wish to optimize
    loss = bce_loss(y_pred, y_target)

    #Step 4: Propagate the loss signal backard
    loss.backward()

    #Step 5: Trigger the optmizer to perform one update
    optimizer.step()
```

# Auxiliary Training concepts

# Auxiliary Training concepts: Evaluation Metrics

❖ **Accuracy**, precision, recall, F1, etc.
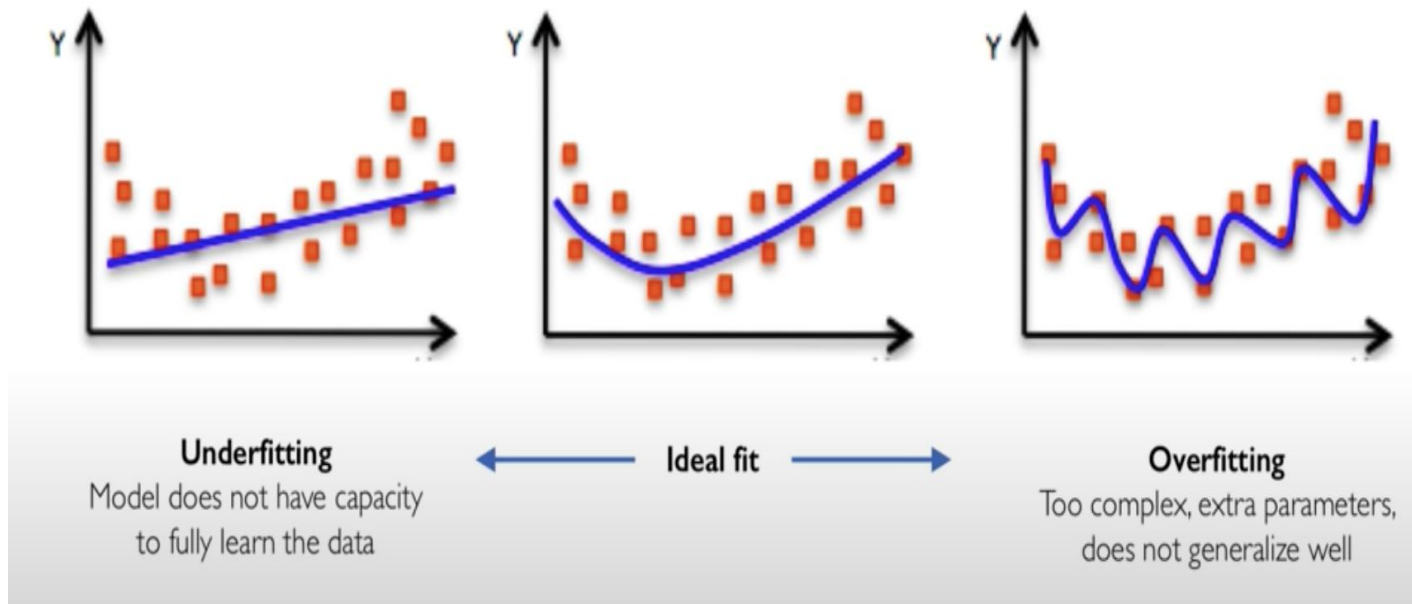
# Auxiliary Training concepts: Splitting the dataset

❖ Standard practice: training, validation and test splitting

❖ K-fold cross validation (for small datasets)

❖ Common split percentage: **70%** training, **15%** validation, and **15%** for testing

❖ For benchmark tasks: predefined training, validation, and test split might exist

➢ E.g. Glue benchmark

# Auxiliary Training concepts: Knowing when to stop training



Stop training before we have a chance to overfit

Loss

Stop training here!

**Legend**

Testing

Training

Training Iterations

From: http://introtodeeplearning.com/slides/6S191_MIT_DeepLearning_L1.pdf

# The problem of overfitting



**Underfitting**
Model does not have capacity
to fully learn the data

Ideal fit

**Overfitting**
Too complex, extra parameters,
does not generalize well
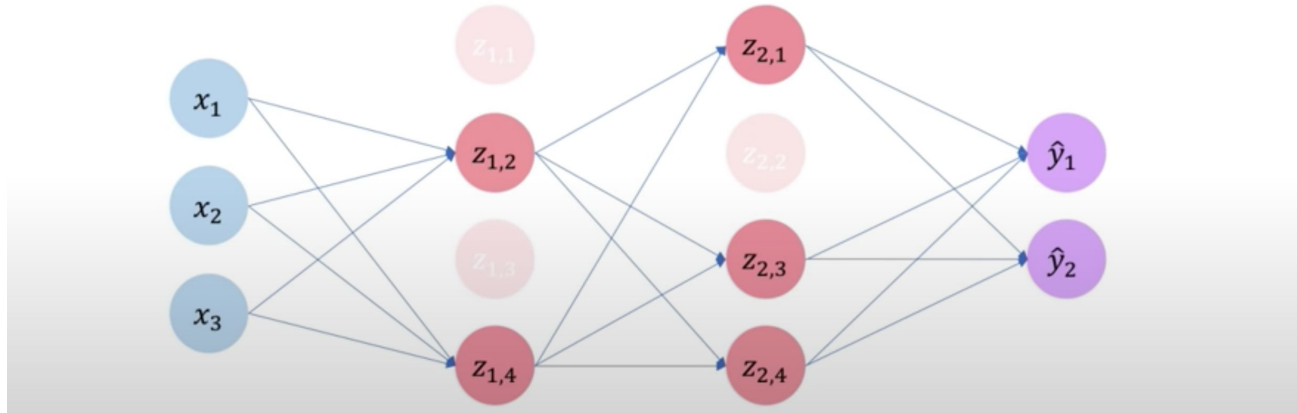
# Auxiliary Training techniques: Regularization

❖ Constrains optimization problem to discourage complex models
❖ Improves generalization on unseen data

# Auxiliary Training concepts: Regularization

❖ Dropout

During training, randomly set some activations to 0
- Typically 'drop' 50% of activations in layer
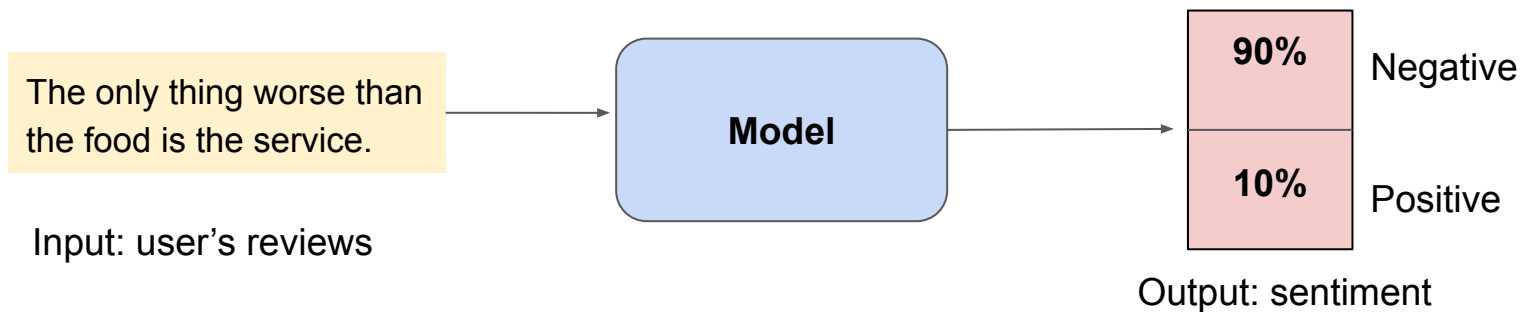- Forces network to not rely on any 1 node



From: http://introtodeeplearning.com/slides/6S191_MIT_DeepLearning_L1.pdf

# Auxiliary Training concepts: Regularization

❖ L2 regularization (weight-decay), L1 regularization
❖ Data augmentation

# Example: Classifying Sentiment of Restaurant Reviews

❖ Task: Classify reviews of a restaurant
  ➢ Positive or negative

The only thing worse than the food is the service.

Input: user's reviews

**Model**

| 90% | Negative |
| 10% | Positive |

Output: sentiment

Yelp dataset simplified version : http://xzh.me/docs/charconvnet.pdf

# Example: Classifying Sentiment of Restaurant Reviews

⇒ Data preprocessing notebook

⇒ Classifying reviews notebook